

Institute for Biological Physics

University of Cologne

Master Thesis

The Fitness Landscape of Translation



Written by

Mario Josupeit

Supervisor: Prof. Joachim Krug

Second revisor: Prof. Andreas Schadschneider

Cologne, 07.08.2020

Contents

1. Introduction	1
2. Biological concepts	2
2.1. Redundant amino acids	2
2.2. Mutations and synonymous mutations	4
2.3. Fitness and fitness landscapes	5
2.3.1. Neutral mutations	7
2.4. Epistasis	7
2.4.1. Sign epistasis	8
2.4.2. Non-epistatic case	9
3. TASEP	11
3.1. General TASEP mechanics	11
3.2. Measures in the TASEP system	12
3.2.1. Density	12
3.2.2. Current	12
3.2.3. Travel time	13
3.3. Homogeneous TASEP	14
3.4. Three phases of the homogeneous TASEP	14
3.4.1. Phase boundaries	15
3.4.2. Edge effects	16
3.5. Inhomogeneous TASEP	17
3.5.1. TASEP systems with random jump rates	18
3.6. Bottlenecks	18
3.6.1. Measures in the TASEP with bottlenecks	19
3.6.2. Phase space of the travel time of a system with a bottleneck	21
3.6.3. Interacting bottlenecks	23
3.6.4. Relation between average density and bottleneck rate	25
4. (Non-)Monotonic parameters in a TASEP with bottlenecks	26
4.1. The current is monotonic in the jump rates	26
4.2. Small fully random systems analyzed with a power series approximation and numerical simulations	28
4.2.1. Calculating sign epistasis numerically	30

4.3. Interfaces in the phase space where travel times are equal	31
4.3.1. Interface where both bottlenecks have the same effect on travel time $\bar{\tau}_1 = \bar{\tau}_2$	31
4.3.2. Interfaces where a bottleneck has the same travel time as the homo- geneous case	32
4.4. Comparisons of phase space interfaces to numerical data	32
5. Modelling interacting bottlenecks	35
5.1. Types of two-dimensional subcubes	39
5.1.1. Inhomogeneous case	40
6. Analysis of the Zwart et al. landscape	42
6.1. Do the assumptions apply to experimental data?	42
6.2. Applying my model	43
6.3. Additive mutational effects within subcubes	46
6.4. Comparing my adapted model to the experimental fitness landscape	47
7. Conclusion and discussion of results	50
Bibliography	52
A. Appendix: Description of code simulating the TASEP	A 1
B. Appendix: Results of the search algorithm for different densities after the bottlenecks	B 3

1. Introduction

In this thesis I examine the fitness effects in the translation step of protein synthesis. The idea for this topic originates from the surprising findings of Zwart et al. in 2018 [35]. Their paper on the *TEM-1* β -lactamase gene of the *Escherichia coli* bacterium states, that synonymous mutations, which are those mutations, that change the nucleic acids, but leave the encoded protein the same, can have a strong fitness effect, with the fitness being the number of offspring per individual. The fitness in an environment with the antibiotic *cefotaxime*, was measured for all combinations of 4 synonymous mutations in the *TEM-1* gene. This fitness is the antibiotic stress resistance $IC_{99.99}$. The synonymous mutations observed are at the 9th, 17th, 87th and 89th codon of the gene that has a length of 284 codons. The space of all 2^4 possible combinations of mutations is called a fitness landscape, connecting a point within the mutations landscape to the measured fitness.

The fitness landscape of synonymous mutations from Zwart et al. [35] features many neutral mutations, which do not change the fitness, as well as sign epistasis, a feature of the landscape where the effect of a mutation has a different sign on different backgrounds. The key to analyzing and understanding such a landscape, beyond looking at the fitness values themselves, is to examine interactions of mutations in the landscape. This work presents a tool for analyzing these landscapes which could lead to a deeper understanding of the characteristics of synonymous mutations. The goal of this thesis is to formulate a model for interacting mutations and analyze the landscape that inspired this investigation. The road to this goal reaches from the biological basics and the TASEP, a non-equilibrium physics model of translation, via the description of a model proposed by this thesis and comparisons to numerical results and literature, to an analysis of experimental results for a fitness landscape of synonymous mutations. The methods of this thesis reach from analytic approaches to numerical simulations and data analysis.

2. Biological concepts

All living cells use proteins, which makes protein production one of the most basic and essential parts of life. Proteins are the tools of the cell and can have various shapes, sizes and functions. Proteins are long chains of amino acids and are in most cases constructed by the cell itself. For each protein produced by the cell, there is a gene in the genetic code that determines the composition of the amino acid chain. The composition determines the function of the resulting protein. The genetic information is stored as triplets of nucleic acids, the codons. The process of protein production is called the *central dogma of molecular biology*, a term coined by Crick in 1970 [2]. This "dogma" states that the genetic information of the gene, stored as part of the cell's DNA (*deoxyribonucleic acid*) is transcribed into a mRNA (*messenger ribonucleic acid*) and then translated by a ribosome into an amino acid chain, which then folds into the functional protein. One mRNA strand can have many protein producing ribosomes on it at the same time, but they can only move in one direction without overlapping and need to read all information of the gene before the protein is completed. This step of protein production has a large amount of processes associated with it and much of the cells functions are designed to enable this process.

Translation, the construction of new proteins by ribosomes, has three steps.

- 1) **Initiation:** A ribosome attaches to the mRNA and starts the translation by adding the first amino acid to the chain that will become the protein in the end.
- 2) **Elongation:** The ribosome moves along the mRNA and attaches one new amino acid to the amino acid chain for each codon.
- 3) **Termination:** Upon reaching one of three stop codons, the ribosome detaches from the mRNA and the amino acid chain detaches as a new protein from the ribosome.

2.1. Redundant amino acids

Amino acids are attached to the end of the produced protein, according to the sequence of codons on the mRNA. Each codon is translated into exactly one amino acid, or is a stop codon. The stop codon marks the end of the gene and therefore signals that the protein is fully produced.

Each codon consists of three nucleic acids, leading to $4^3 = 64$ specific codons, cf. figure 2.1. Many of these possibilities are redundant, meaning that multiple codons code for the same amino acid. This redundancy is the reason why up to six different codons encode each of the 21 canonical amino acids. Codons that code for the same amino acid are called synonymous

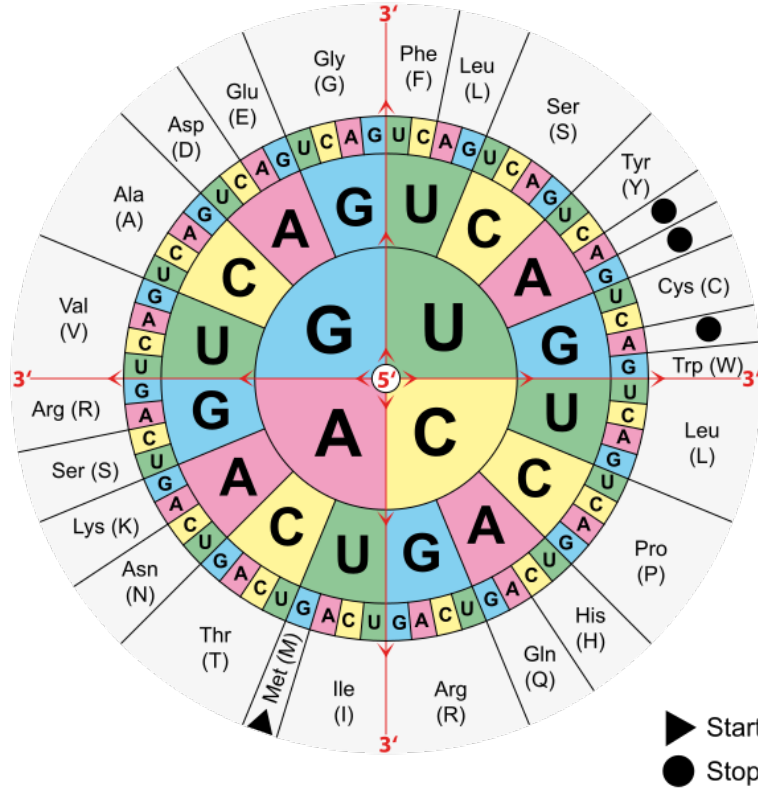


Figure 2.1.: The mRNA codons with their respective amino acids. Note that only the start codon methionine (labeled *Met* (*M*) in the figure) and tryptophane (labeled *Trp*(*W*)) are encoded by a single codon. The other amino acids have synonymous codons which encode them.

Graphic source https://de.m.wikipedia.org/wiki/Datei:Aminoacids_table.svg, visited 25th of July 2020, licensed as free to use (Public domain).

[22]. When first discovered, it was believed that synonymous codons have no effect on protein production, because they do not alter the sequence of amino acids in the protein. More recent experimental findings [35] show a large change in the features of cells after exchanging codons with their synonymous counterparts. The synonymous substitution, i.e. the exchange of one synonymous codon for another, keeps the sequence of amino acids the same. There exist examples for different speeds at which different synonymous codons are read during elongation [31, 18]. This different translation speed can lead to changes in protein stability and production yield [21]. Exploring this phenomenon is very interesting since this is a new angle from which researchers can understand protein production.

2.2. Mutations and synonymous mutations

Mutations are exchanges, deletions or insertions of genetic material in the genetic code of an organism. This means that a mutation changes the genetic code and sometimes a feature of the cell. The effects of mutations can be various. The protein encoded by the gene may change its structure and stop working or, in a rare, but better case for the cell, the protein gains a new feature that is beneficial for the cell. The new feature caused by a mutation can be anything from a metabolic function to the disabling of an important protein leading to the cell's death. If there is a measurable change, the phenotype of the cell changed. The concept of phenotypes is explained in very large detail by Taylor in chapter 6 of [30]. It is important to note, that the effects of a mutation may be beneficial or disadvantageous. Especially if there is no competition between individuals, for example because the individuals never interact with another, less efficient individuals also grow well and produce offspring. For experiments on bacteria, the solution is often diluted so much, that all individual bacterial colonies are the offspring of one ancestor each.

In the following sections I often use the term *sign of a mutation*, where the sign is positive, if the mutation is beneficial, and negative if it is disadvantageous. Difference in genetic code is used as a measure for distance between *genetic species*. The more nucleotides differ from one genome to the other, the further they are away from another. the members of one genetic species are all individuals in a population, that share the exact same genetic code [1], which is a more rigid definition of the term *species*. This also means that mutations change one genetic species into another. The term *species* is only defined for sexually reproducing organisms and not to be confused with the term I use in this work.

In this thesis, I focus on synonymous mutations. This is a special case, where codons are changed in such a way, that the amino acids they code for stay the same. Even though there is no change in the protein sequence, there are interesting effects on the phenotype [22]. This is an example for effects on the phenotype of an organism beyond changing its amino acid sequences. A feature that changes when a synonymous mutation happens is the elongation rate of the changed codon. The elongation rate is the rate at which the ribosome translates a codon and is an essential factor for the organisms fitness, which is explained in the next section.

2.3. Fitness and fitness landscapes

Fitness is a macroscopic variable dependent on the genome of the individual. The genome is the collection of all information encoded by the genes of an organism. It gives a single macroscopic value dependent on many microscopic values, similar to the free energy of a gas or the color of a crystal. It is described as a function of the full genome and the living circumstances of the organism. It is impossible to define a fitness without knowing the environment of the organism. In the most simple case, one can find certain proteins that are the most important for the survival of the cell. In the example study referenced multiple times in this thesis by Zwart et al. [35] the *E. coli* bacteria grow in a medium with a high antibiotic concentration, which is why they need to produce a certain protein, which deactivates the antibiotic, or die. The fitness $F(\vec{\nu})$ of this organisms then is highly dependent on the codon sequence $\vec{\nu}_0 = (\nu_1, \nu_2, \dots, \nu_L)$ of length L of the gene that can deactivate the antibiotic. The fitness is then $F(\vec{\nu}_0) = F_0$. A mutation m_1 changes the nucleic acid sequence at position x_1 . This could be a point mutation exchanging only one nucleic acid, or a insertion or deletion of a section in the gene. The important sequence change for this thesis is the single exchange of one amino acid, which is also the most common mutation, which is why in the following only one rate and one position is taken into account at a time.

$$m_1 : \vec{\nu}_0 \mapsto \vec{\nu}_1 \quad (2.1)$$

The new genome is then $\vec{\nu}_1 = (\nu_1, \nu_2, \dots, \nu'_{x_1}, \nu_{x_1+1}, \dots, \nu_L)$ and may exhibit a different fitness

$$F(\vec{\nu}_0) \neq F(\vec{\nu}_1) . \quad (2.2)$$

A visualization of the fitness values and also the genomic distances is the fitness landscape. It consists of different genomes mapped to their associated fitnesses. Two points in this space represent two different genetic species and are connected via mutations that change one genetic species into the other. Because each visualization shows different features of a fitness landscape, there are many ways of visualizing it.

Some common visualizations are the fitness plotted against the number of mutations and the N-dimensional hypercube.

The fitness plotted against the number of mutations shows the distance from the original genetic species. The original species is called the *wildtype*. The distance is the number of mutation steps that need to be taken to get from the wildtype to the other mutants. It is

2.3 Fitness and fitness landscapes

most commonly used to show that a landscape is very smooth or very rough, because this feature of the landscape is visualized very well. An example for fitness plotted against the number of mutations can be seen in figure 2.4.

The fitness landscape can also take the form of an N-dimensional hypercube, where N is the number of different mutations of the organism [6]. This structure spans an N-dimensional space of edge length 1 (so each site can either be mutated or not mutated), with the 2^N fitness values at the corners of the hypercube. This visualization is most commonly used to show the pathways along which mutations can get from one point in the landscape to another. The fitness values themselves are less prominent in this visualization. Many examples of this visualization can be found in chapter 5 and in the example 2.2. The general form of this hypercube has four values on each edge. These four values correspond to the four nucleic acids that can be at those edges. It is common practice though to use a binary alphabet if there is only a maximum of one mutation at each position on the gene observed. The fitness landscape is easier to display with only two points on each edge.

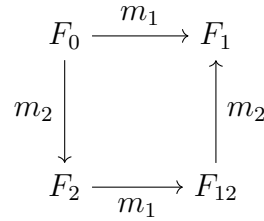


Figure 2.2.: A two-dimensional fitness landscape displayed as a hypercube. Each direction in the landscape is one mutation. In the upper left corner is the unmutated wildtype fitness, that is the fitness of the original genetic species. Arrows point to the nodes of higher fitness. As an demonstration example I choose $F_0 < F_2 < F_{12} < F_1$.

Both visualizations can show how mutations interact, which is a main focus of this thesis and explained in the next section on epistasis.

There are many parameters that can be treated as the fitness of an organism. Whether it is the resistance to an antibiotic, if the cell grows in a medium with antibiotics, the ability to process more nutrients, if there is a new source of energy available, or the time it takes for an individual to produce offspring. All of those measures lead to reproductive success and can be a proxy for fitness. That is why, in evolution, fitness means reproductive fitness according to Wright [33]. This is the amount of offspring of one particular genetic species in the next generation, which is the definition of fitness used in this thesis.

2.3.1. Neutral mutations

If a mutation has no effect on fitness, it is called neutral (or silent). Displayed similarly to figure 2.2, figure 2.3 shows two cases of neutral mutations. In figure 2.3a, the mutation m_2 is always neutral. In figure 2.3b, the mutation m_2 is only neutral, if the mutation m_1 happened before it. Both cases do exist in nature. Displayed as a fitness landscape with the fitness plotted against the number of mutations, neutral mutations are lines without a slope, also called "flat". I use this to describe mutations in chapter 5. The case from figure 2.3b is further discussed in the next section 2.4, which explains epistasis.

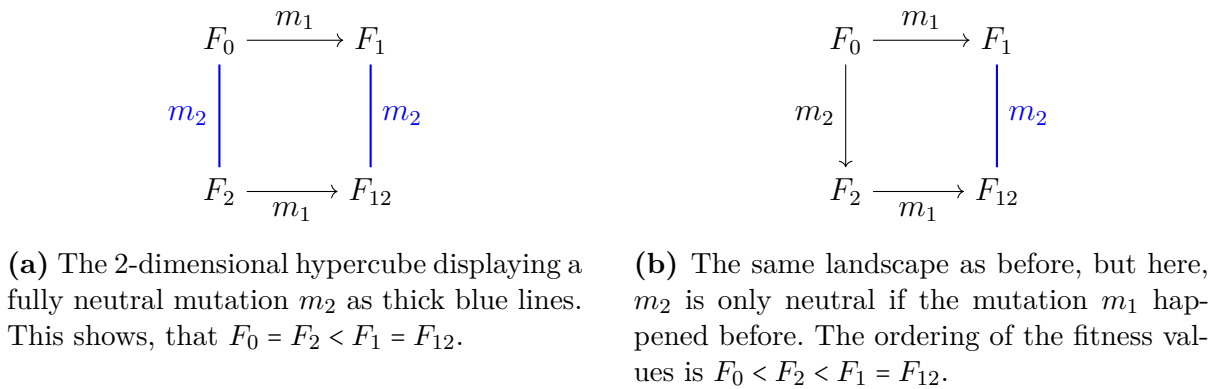


Figure 2.3.

2.4. Epistasis

The effect of mutations often depend on the genetic background, i.e. all other information encoded by the genome of the organism. A very intuitively accessible example is a bacterium that develops both the ability to gather a new food and separately the ability to digest it. This bacterium could develop both traits independent from another, but the large benefit only exists if both traits are present in the same bacterium at the same time.

This very important concept for this thesis is called epistasis. It describes the interaction of mutations. There are many different definitions for epistasis depending on the aspect of interest. The definition that I use states that epistasis is the change on the effect of one mutation m_2 due to the presence of another mutation m_1 , described on page 667 in Crow and Dove's book [3]. This definition can also be applied for interactions between more than two mutations. In the previous section there are some examples for (non-)epistatic landscapes. Figure 2.3a shows no epistasis, figure 2.3b shows epistasis in the mutation m_1 , because the effect of m_1 is either neutral or beneficial, depending on the presence or

2.4 Epistasis

absence of mutation m_2 and figure 2.2 shows the special case of sign epistasis, where the presence of mutation m_1 changes the sign of the effect of mutation m_2 compared to the case where mutation m_2 acts on the wildtype. The definition for the terms sign epistasis and non-epistatic follow in the subsections 2.4.1 and 2.4.2.

2.4.1. Sign epistasis

A special case of epistasis is sign epistasis. It describes the change of one sign of a mutation dependent on the presence or absence of another mutation. This non-monotonic effect is very interesting, because two positive effects may produce a negative effect in conjunction or two negative effects may be beneficial together. Weinreich explains, that the sign of a mutation is under epistatic control [32].

If F_0 is the fitness of the unmutated wildtype species with genome $\vec{\nu}_0$, F_1 is the fitness after mutation m_1 on the genome $\vec{\nu}_0$ took place, turning it into $\vec{\nu}_1$. F_2 is the fitness after mutation m_2 emerged on $\vec{\nu}_0$ and F_{12} is the fitness with both mutations present. The effect of mutation m_1 on the background $\vec{\nu}_0$ is equal to the fitness difference $F_1 - F_0$. The effect of m_1 on $\vec{\nu}_2$ is equal to the fitness difference $F_{12} - F_2$.

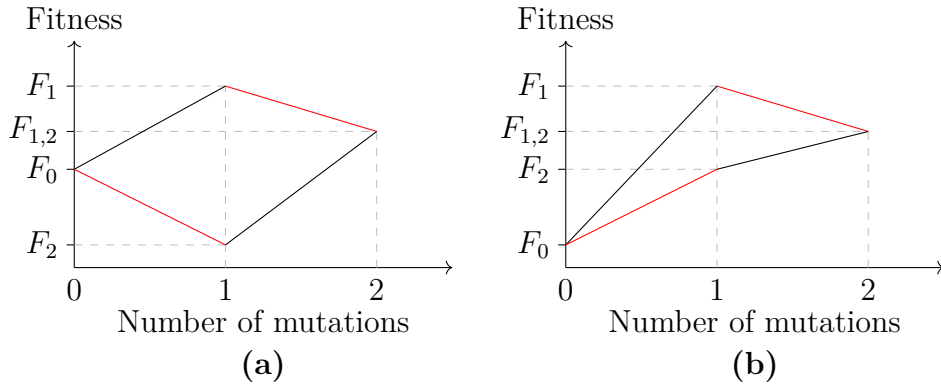


Figure 2.4.: The epistatic effects of two interacting mutations. **(a):** An example landscape with monotonic fitness effects, i.e. no sign epistasis. Mutation m_2 (red lines) always has a negative effect and m_1 always has a positive effect on fitness. The angles of the mutations change slightly, depending on the presence of the other mutation, but the sign does not change. **(b):** Non-monotonic system exhibiting sign epistasis. m_2 (red lines) has a positive effect only if m_1 is not present. m_1 itself does not display sign epistasis with respect to m_2 because it always has a positive effect whether m_2 is present or not.

If the difference between the fitness values without the mutation m_2 being present ($F_1 - F_0$) and with m_2 , ($F_{12} - F_2$), have different signs, the presence of the mutation m_2 changes the

sign of the effect of the mutation m_1 . Therefore if either or both

$$\begin{aligned} E_1 : (F_{12} - F_1)(F_2 - F_0) &< 0 \\ E_2 : (F_{12} - F_2)(F_1 - F_0) &< 0 \end{aligned} \tag{2.3}$$

are true, the presence of one mutation flips the sign of the effect of the other.

Because the definition allows for either one of the equations or both to be true, drawn as a hypercube this sign epistasis shows up as antiparallel arrows. In the example 2.2, $(F_{12} - F_1)(F_2 - F_0) < 0$ is true and therefore the arrows associated with the m_2 mutation are antiparallel. In the example 2.4b, the mutation m_2 has a positive effect, when it emerges on the wildtype ($F_2 > F_0$), but has a negative effect when the mutation m_1 is already present ($F_{12} < F_1$). This can be seen in the different sign of the slopes of the red lines in example 2.4b. The example 2.4a does not feature this effect.

Sign epistasis is very interesting since one mutation can have a very large impact on the effect of another mutation. It not only changes the effects strength, but even whether the effects is positive or negative for the organism and the effects are not constant.

For this thesis sign epistasis is important, since the experimental fitness landscape by Zwart et al. [35] displays many cases of sign epistasis.

2.4.2. Non-epistatic case

In the non-epistatic case mutations are independent. Therefore the fitness effects are additive. In contrast to the epistatic case, the fitness effects of each mutation is constant. If

$$F_{12} - F_2 = F_1 - F_0 \tag{2.4}$$

is true, there is no epistasis between m_2 and m_1 . Transforming equation 2.4, it is equivalent to $F_{12} - F_1 = F_2 - F_0$. Therefore both mutations have a constant effect on fitness and are independent from each other. Depicted as a hypercube of two dimensions this is shown in figure 2.5. The signs of the mutational effects C_1 and C_2 are not important for equation 2.4 to hold, therefore the strength of the effect remains constant. The effects C_1, C_2 can have positive or negative signs.

2.4 Epistasis

$$\begin{array}{ccc}
 F_0 & \xrightarrow{m_1} & F_1 = F_0 + C_1 \\
 m_2 \downarrow & & \downarrow m_2 \\
 F_2 = F_0 + C_2 & \xrightarrow{m_1} & F_{12} = F_0 + C_1 + C_2
 \end{array}$$

Figure 2.5.: A fitness landscape with constant mutational effects C_1, C_2 .

3. TASEP

In this thesis, I simulate ribosomes movement on the *mRNA* with the **t**otally **a**symmetric **s**imple **e**xclusion **p**rocess (TASEP) model, which is a well established model for protein synthesis and a standard model in the area of non-equilibrium physics. It was suggested in 1968 by MacDonald and Gibbs [20] and is also used for simulating other transport processes like traffic jams [23] and myosin movement (e.g. [8, 12, 15, 20, 19, 26]). Even though the kinetics of the TASEP are simple, it displays very interesting effects, such as spontaneous shock formation, phase transitions and edge effects. The phases of a TASEP and the fluctuations that occur within them during simulations are the topic of a paper by de Gier and Essler [9] in the context of solid state physics. For traffic models, all of these effects are easily observed in real scenarios, for the biological process that sparked the idea, there are many obstacles to the observation because most ways of measuring the parameters of the motion in the system require the system to be stopped from working. Knowledge about the macroscopic observables of the TASEP system are important to gain an understanding of the mechanics of it, these are explained in section 3.2.

3.1. General TASEP mechanics

The TASEP exists on a one-dimensional lattice of length L on which particles move unidirectionally. The particles are subject to hardcore interaction, so they can not overlap or overtake another. They move from left to right taking steps, respectively jumps, of distance 1 along the one-dimensional system. These jumps occur between the sites and reflect the one elongation step of the ribosome. After entering at the left end with jump rate α they traverse the system at rates $(\omega_1, \omega_2, \dots, \omega_{L-1}) = \vec{\omega}$, which are the $L - 1$ jump rates between all L sites in the system. The particles exit on the right, from site L , with jump rate β . The jump rates determine particle movement, given there is an empty space to the right of them. At its starting position the system is connected to an infinite reservoir of particles, so there is always a particle able to fill the first site. At the exit it is connected to a particle sink, so a particle at the last site can always leave the system [34].

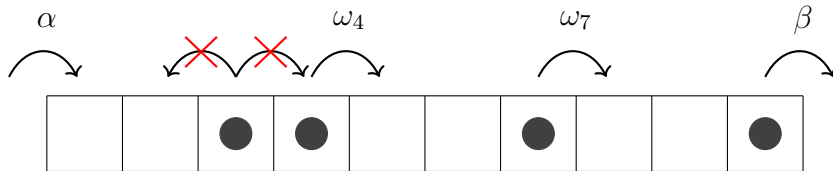


Figure 3.1.: Schematic representation of the allowed movement in a TASEP of length 10. Particles enter at rate α on the left, move at their local rate and can not jump backwards or occupy the same spot as another particle. They leave the system on the right at rate β .

3.2 Measures in the TASEP system

3.2. Measures in the TASEP system

For the TASEP the interesting macroscopic parameters are the stationary current J , the stationary average density $\bar{\rho}$ and the travel time $\bar{\tau}$. In the following I describe their general formulation. There is an exact solution for the homogeneous case in section 3.3. These parameters change when bottlenecks are present, explained in section 3.6.1. There is also an approximation for the current and density for random systems with a low initiation rate α by Szavitz-Nossan et al. [28] which is the topic of section 4.2.

3.2.1. Density

The local density ρ_i in the steady state, i.e. the system after it has relaxed for a sufficient time, is the likelihood to find a particle at position i . Its average across the system

$$\bar{\rho} = \sum_{i=1}^L \frac{\rho_i}{L} \quad (3.1)$$

is a measure for the average amount of particles in the system. Even though it has been used as a measure for fitness, this density does not measure fitness. This is mostly used in experiments where ribosome profiling, which is a method where ribosomes are used to shield *mRNA*, is performed to measure the density of ribosomes. This misconception is often based on the idea that ribosomes move simultaneously, or the assumption that translation is only limited by α , which is a setup that has, according to my knowledge, not been observed in an biological system. If the density of ribosomes is high, the cell has less ribosomes available, costing energy, and the protein production yield does not increase as described in more details by Plotkin and Kudla in the section on measurements in their paper [22].

Closely related to the density is the average hole density, which is the likelihood to not find a particle at position i , which is simply $1 - \rho_i$.

3.2.2. Current

The current J in a TASEP is generally a function of all rates ω_i . There is no general solution for it, but at any moment in time, it can be understood as the rate at which the average particle in the system moves. This definition does not give an analytic formula for the current though, because in an inhomogeneous system, the rates at which the particle leave their sites is highly variable. In large homogeneous systems it can be approximated

as the product of the average density with the average hole density,

$$J = \bar{\rho}(1 - \bar{\rho}) . \quad (3.2)$$

The current can be approximated by the maximal permitted current $J_{\omega_{\min}}$ through the site with the lowest rate ω_{\min} .

In the case of protein production, the average current is a measure of how much protein is produced per *mRNA* strand per unit time and therefore an often used measure for fitness.

3.2.3. Travel time

The travel time T as formulated by Szavits-Nossan and Evans [29], is a measure for the time it takes a particle from the first position in the TASEP to leaving the system. The travel time is the sum of the local densities of particles $\bar{\rho}$ divided by the current J . This thesis uses a slightly different formulation, because I approach the topic of travel time from another direction, but the definitions are equivalent. The difference from the one in the paper by Szavits-Nossan and Evans is due to them starting to count at the second site, while I start at the first, which is why I multiply by L instead of $L - 1$ and that I focus in the average travel time per site $\bar{\tau}$. The travel time T is

$$T = \frac{L\bar{\rho}}{J} . \quad (3.3)$$

The time that a ribosome spends at each site i is

$$\tau_i = \frac{\rho_i}{J} , \quad (3.4)$$

and the average time that a ribosome spends at a site is

$$\bar{\tau} = \frac{\bar{\rho}}{J} . \quad (3.5)$$

In the context of translation, the travel time is a measure of the time it takes one ribosome to produce one protein. Each particle encounters other particles with the probability $\bar{\rho}$. So if there is more jamming in the system, then the travel time is longer. The inverse of the travel time $\frac{1}{\bar{\tau}}$ is the translation efficiency. It is the rate at which proteins are produced per ribosome and is another measure for the fitness of an organism.

3.4 Three phases of the homogeneous TASEP

3.3. Homogeneous TASEP

The TASEP is called homogeneous if the jump rates for all sites are homogeneous,

$$\omega_1 = \omega_2 = \dots = \omega_{L-1} =: \omega . \quad (3.6)$$

There are three free parameters in all homogeneous TASEP systems, the initiation rate α , termination rate β and the homogeneous elongation rate ω . The rates are only relative to an arbitrary time measure and the results are the same after renormalizing with $\alpha^* = \frac{\alpha}{\omega}$ and $\beta^* = \frac{\beta}{\omega}$, which is why any homogeneous ω can be set to 1 after rescaling. The homogeneous TASEP has been analytically solved by Derrida et al. and the stationary current and density are known [4, 5]. The following section 3.4 sums up the results for the homogeneous TASEP.

3.4. Three phases of the homogeneous TASEP

The homogeneous TASEP system separates into three phases, the low density (LD), high density (HD) and maximum current (MC) phase. The phase transition between the high density and low density phase is called the shock phase (cf. figure 3.2a).

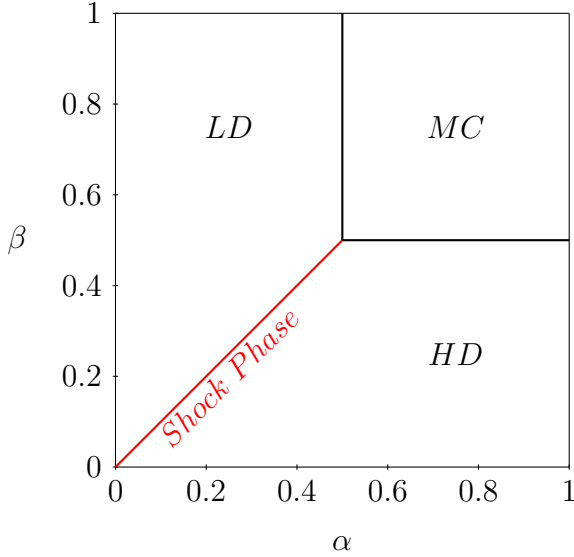
If α is smaller than 0.5 and β is larger than α , the system is in the low density phase. α is the rate limiting factor for ribosome movement. For this case, the termination rate β is not a current limiting factor, because the density ρ in the system is always lower than the rate at which ribosomes exit from the last site. The rates in the bulk do not limit the movement because they allow $J_{MC} = 0.25$, the rate of the last site allows $J_{\text{last site}} = \beta(1-\beta)$ and the current through the start site is $\alpha(1-\alpha)$, which is smaller than the other two. Therefore the density in the system is α . This leads to a lower density than in all other systems, hence the name **low-density** phase (LD).

If β is smaller than 0.5 and α is larger than β , the system has a very high density. In contrast to the low-density phase the current is limited by the termination rate β . The density in this system is $1-\beta$, because particles leave the system at rate β , and the density that remains at the last site and the traffic jam propagates to the left is $1-\beta$. The system supplies more particles than can exit and is in a phase of high density (HD).

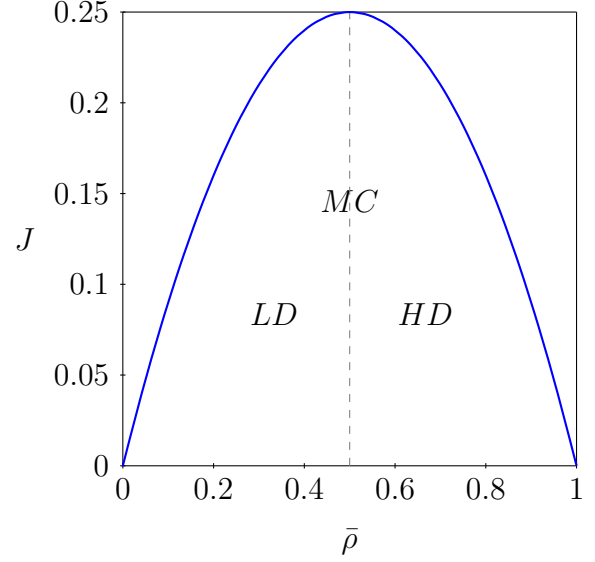
The third phase is the **maximum current** phase (MC). If a sufficient amount of particles can enter and leave the system. The entry and exit no longer limit the travel in the system. This is true if the entry rate $\alpha \geq 0.5$ and the termination rate $\beta \geq 0.5$. In this case, the

3.4 Three phases of the homogeneous TASEP

current in the bulk is now rate limiting, because it reaches its maximum of 0.25 (cf. figure 3.2b). The density $\bar{\rho}$ is 0.5 in the whole phase.



(a) The phase diagram for the homogeneous TASEP. It splits into three phases, high density (*HD*, lower right), low density (*LD*, upper left) and maximum current (*MC*, upper right). The boundary between high density and low density is the shock phase (red), where the two bordering regions coexist.



(b) The connection between the current J and the average density $\bar{\rho}$. On the left of the peak is the LD phase, on the right the HD phase. The line in the middle signifies the MC phase.

3.4.1. Phase boundaries

At the intersection between the phases, phase transitions occur. Between the low-density and the maximum-current phase and between the high-density and the maximum-current phase, there are second order phase transitions. The low-density system fills with particles as α increases until it reaches $\alpha = 0.5$, where the bulk can no longer support a higher current and the maximum current phase is reached. The opposite is true for the high-density phase, here the density decreases as β increases until it becomes a maximum current system at $\beta = 0.5$.

The phase transition between the low density phase and the high density phase is different. Here the system approaches the line $\alpha = \beta$ from either the low-density or high-density phase. There it enters the shock phase, in which the system splits into a low-density part at the start and a high-density part at the end. The two phases coexist, because the particle enter at the same rate as they leave, so the system neither fills up nor drains. The particles enter

3.4 Three phases of the homogeneous TASEP

at a low rate and have almost no other particles in their way, due to the low density in the first part of the system, so they reach the intersection between the two parts rather fast. At the other end of the system, the particles leave at a low rate, leading to a traffic jam in front of the termination end. As soon as there is a vacant spot at the last site, due to the high density at the end, this new hole is transported to the right very fast. The intersection between the phases is called a shock, due to the sudden change in density.

The shock diffuses through the system. Whenever a new particle arrives at the shock, the shock moves towards the start of the system, whenever a hole reaches the shock, it moves towards the end. This diffusion leads to a towards the end of the system linearly increasing average density ρ_i (cf. figure 3.3).

3.4.2. Edge effects

In general TASEP systems there are always edge effects. The density at the borders decays into the system. If the system is in the maximum current phase, this decay is a power law. If the system is in the low-density or high-density phase, it decays exponentially. These changes in density is completely relaxed, there are tails at the boundaries like in figure.

In the homogeneous TASEP, the larger the initiation rate α is, the stronger are the edge effects close to the start site. The sites close to the start display a density that is larger than the density in the center of the system if $\alpha \geq \beta$. A similar effect can be observed at the end of the system for $\beta \geq \alpha$, where the density drops. This is visualized in 3.3.

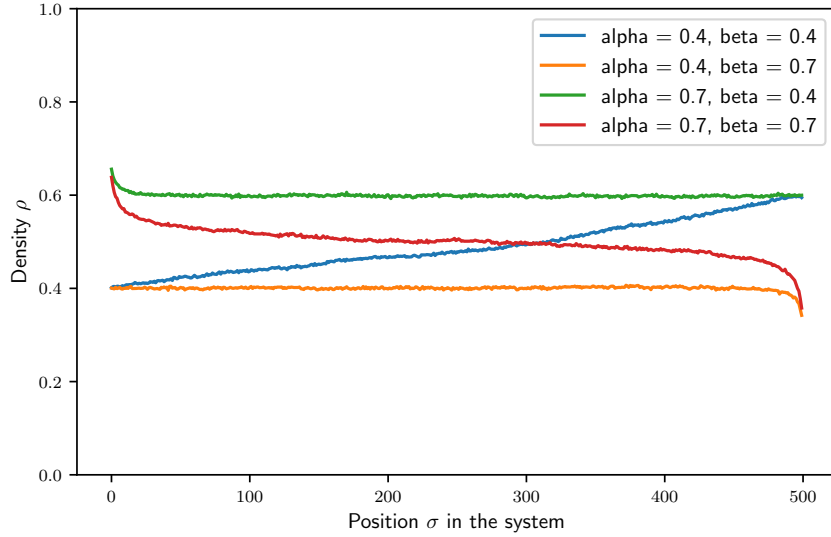


Figure 3.3.: Example runs for the homogeneous system. The graphic shows the distinct characteristics of each phase. In the shock phase (blue) the density is monotonically increasing. In the high/low density phase (green/orange), edge effects are visible at the start/end and the system has the same overall density otherwise. In the maximum current phase (red), the density is on average at 0.5 and has tails at both ends.

Edge effects exist at all rates bordering different rates within the system as well and depends on the difference between the rates. Therefore this effect can also be observed in the bulk of inhomogeneous systems, because not all rates ω_i are the same. Other features of the inhomogeneous TASEP are explained in the next section 3.5.

3.5. Inhomogeneous TASEP

A synonymous mutation can change the elongation time of the affected codon and therefore the rate at that site. From experiments it is known that the change in the rates due to synonymous mutation can differ by a factor of up to 4 [31, 18]. These different rates have to be reflected in the simulations. The inhomogeneous TASEP reflects, that the rates ω_i of different sites i , can have different values.

Unfortunately, in contrast to the homogeneous TASEP, these systems are not solved. The only analytical solution is an approximation for systems that have one very small rate α or ω_i by Szavits-Nossan et al. [28]. Therefore numerical simulations are required to understand these systems. Even for relatively small finite systems, there is no general solution and

3.6 Bottlenecks

numerical solutions are required to understand their behavior.

3.5.1. TASEP systems with random jump rates

The first approach to simulating a system of many free parameters is generating a system with random rates. These systems have large statistical noises, which do not abate, even at long timescales. The causes for the noise are edge effects throughout the system, that exist whenever different rates border another.

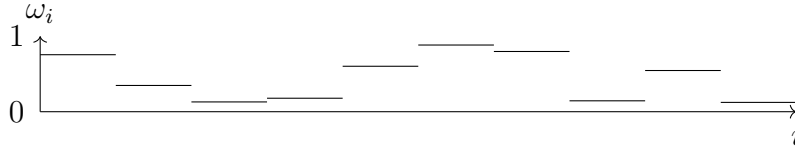


Figure 3.4.: An example of a random landscape with random rates ω_i .

3.6. Bottlenecks

A second approach to understanding the TASEP using numerical simulations is to start from the analytically solved homogeneous TASEP with $\alpha = 1$ and $\beta = 1$ and replace one of the jump rates ω_i at position x in the system with a rate r that is smaller than the other rates. This local inhomogeneity is called a bottleneck. The density in systems with a bottleneck acts similar to a fixed shock, but in contrast to the shock phase, it does not diffuse in the system and the density at the start of the system is high, and at the end of the system is low. Properties of bottlenecks are nicely explained by Schadschneider, Chowdhury and Nishinari in chapter 6 of their book on transport systems [23].

In this section, I make statements that are true for large systems ($L \gg 1$). The discontinuity in the rates leads to a fixed density behind the bottleneck and another fixed density before the bottleneck. I only consider large systems in the following, because the parameters bottleneck rate r and density after the bottleneck ρ can be used interchangeably, cf. section 3.6.4. The bottleneck is fully characterized by its position x in the system and rate r and therefore in the limit of large L , it can also be characterized by the density after the bottleneck and the position. The bottleneck rate r is easier to use in numerical simulations, but has no explicit meaning in the context of fitness, because the exact relation between fitness and the rate is unknown and the density has a meaning connected to fitness, but can not be used as an input parameter in simulations. The exact function $\rho(r)$ is not known for finite systems, but can be numerically calculated. This is done in section 3.6.4.

Janowski and Lebowitz approximate the current J and density ρ depending on bottleneck rate r [13] and give an expansion for these values for finite systems [14]. Szavits-Nossan uses a matrix formulation of the transitions to calculate these functions up to the third order in the lowest rate in the system [27].

3.6.1. Measures in the TASEP with bottlenecks

The measures explained in section 3.2 change in a system with bottlenecks. A single bottleneck in an otherwise homogeneous system separates it into two parts, where each are themselves a homogeneous TASEP. The bottleneck reduces the density behind it and increases the density before it. In large systems, local inhomogeneties around the bottleneck in the density profile can be ignored because the density behind the bottleneck is mostly dependent on the rate r of the bottleneck. In the following I assume the system to be large to have parameters that are more sensible for the model.

The average density of the second part of the system depends on the rate r of the bottleneck (cf. figure 3.10). Because all particles have to travel through both parts of the system, the current of particles is

$$J = \rho(1 - \rho) \quad (3.7)$$

everywhere in the system and the density ρ is mainly dependent on the bottleneck rate. Equation (3.7) has two solutions for $\rho \in (0, 0.5)$, therefore the average local density after the bottleneck $\overline{\rho_{\text{after}}}$ needs to be equal to the average local density of holes before the bottleneck $1 - \overline{\rho_{\text{before}}}$ (cf. figure 3.5), the average density of particles before the bottleneck and after the bottleneck add up to 1. For ease of notation I define

$$\rho := \overline{\rho_{\text{after}}} , \quad (3.8)$$

$$\Rightarrow \overline{\rho_{\text{before}}} = 1 - \rho . \quad (3.9)$$

There is a similarity to the shock phase, described in section 3.4.1, because the system is separated into two parts by the bottleneck, one high-density phase and one low-density phase, but the shocks in bottleneck systems do not diffuse through the system like in the shock phase, but are fixed.

The position of the bottleneck in the system is given by

$$x := \frac{i}{L} \Big|_{\omega_i=r} , \quad (3.10)$$

3.6 Bottlenecks

which is a value between 0 and 1. For a bottleneck with rate r at site i in a large system of length L , the average density of the whole system is the length x times the density before the bottleneck $1 - \rho$ plus the length after the bottleneck $1 - x$ times the density after the bottleneck ρ . It is

$$\bar{\rho} = (1 - \rho)x + \rho(1 - x) . \quad (3.11)$$

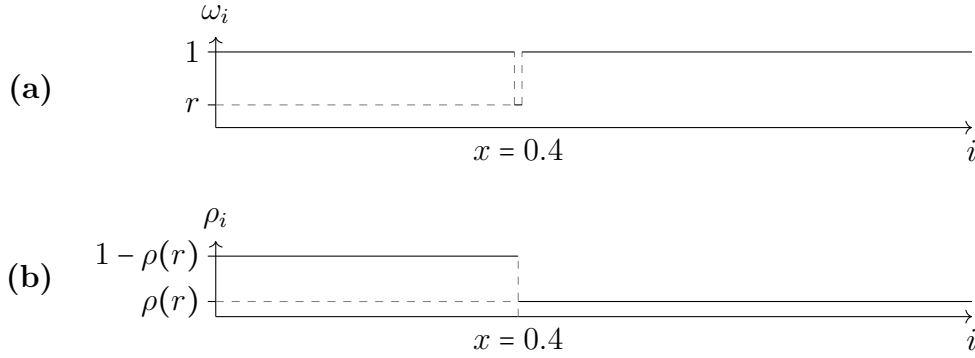


Figure 3.5.: (a): An example for the rates of a TASEP with a bottleneck at position $x = 0.4$ of rate r . (b): Schematic representation of the density profile of the TASEP with a bottleneck. The system is separated into two parts, a high-density system at the start and a low-density system at the end.

The average travel time per site is

$$\bar{\tau} = \frac{\bar{\rho}}{J} = \frac{x}{\rho} + \frac{1 - x}{1 - \rho} \quad (3.12)$$

$$\Leftrightarrow \bar{\tau} = \frac{1}{1 - \rho} + x \left(\frac{1}{\rho} - \frac{1}{1 - \rho} \right) . \quad (3.13)$$

The really interesting feature of the travel time is shown in figure 3.6. The system with high jump rates all throughout the system in figure 3.6e is not the fastest, since it is in the maximum current phase where for every particle the probability to have its path blocked is equal to the average density of the system $\rho = 0.5$. The fastest moving particles move through systems with a bottleneck right at the start in figure 3.6a, that prevents jamming all throughout the system. The lower current due to the bottleneck is overcompensated by the low density in the second part of the system.

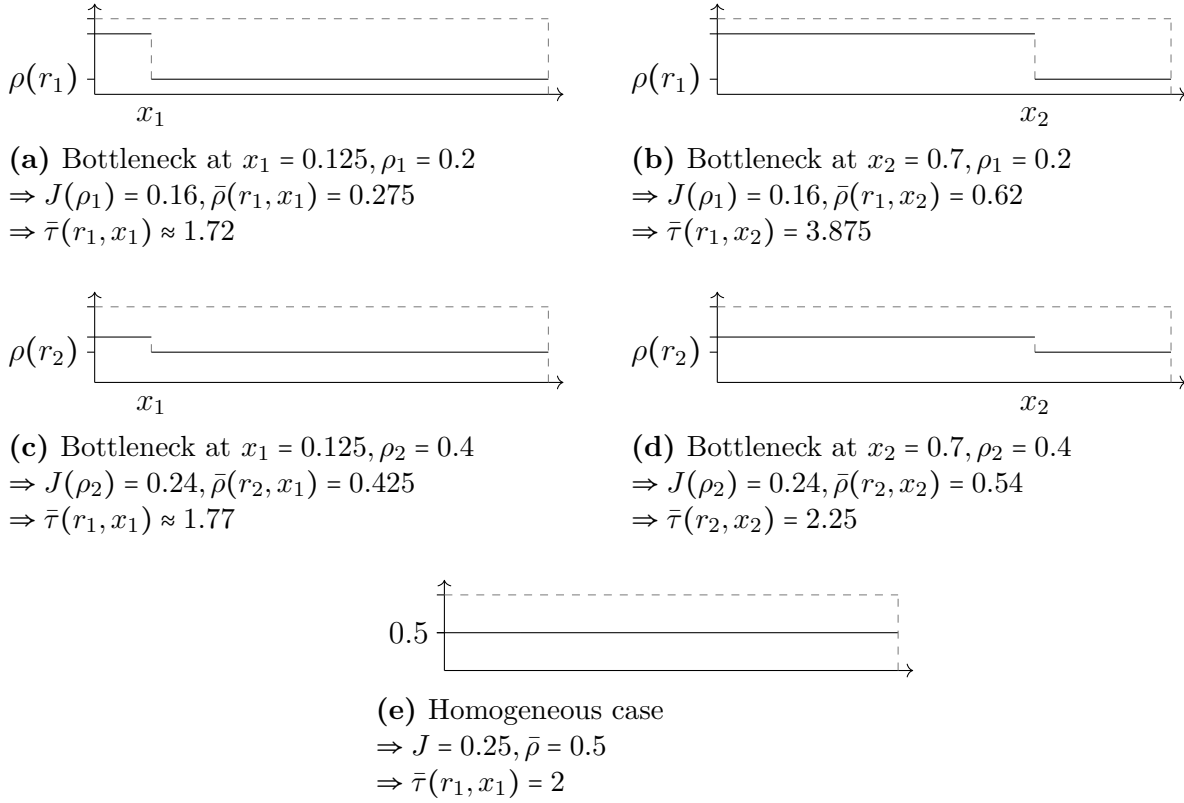


Figure 3.6.: Comparison of average density $\bar{\rho}$, current J and travel time $\bar{\tau}$ of four different bottleneck setups with bottleneck locations x_1, x_2 , densities after the bottleneck ρ_1, ρ_2 (a)-(d) and the homogeneous case (e). The travel time is fastest in (a), because the small rate at the start of the system prevents jamming. The setup (c) is slower because the rate of the bottleneck is higher than in (a), causing a higher density after the bottleneck ρ and therefore increases the likelihood of jamming. If the bottlenecks are at the end of the system, the higher density setup (d) has a smaller travel time than (b).

The travel time is a function of both the location and the rate at the bottleneck. More details are discussed in the next section on the phase space of the travel time dependent on the defining parameters of the bottleneck, x and r .

3.6.2. Phase space of the travel time of a system with a bottleneck

The right edge of the phase space in figure 3.7 is the homogeneous system ($\rho = 0.5$). At this line, it is as if a rate $r = 1$ were inserted into the system, leaving it homogeneous because there is in fact not bottleneck inside of the system. All but the travel time for the case where the lower current due to the insertion of a bottleneck is compensated by the lower average density ($\bar{\tau} = 2$), never reach the line $\rho = 0.5$. For any fixed $\bar{\tau} < 2$, the density and

3.6 Bottlenecks

location of the bottleneck only exist in a certain interval for both ρ and x ¹. If $\bar{\tau} > 2$, the location of the bottleneck can be anywhere in the system, but there still is a maximum density if the bottleneck is at the end of the system.

In the case where the density ρ is fixed to a constant value ρ_c , the travel time is

$$\bar{\tau}_c = \frac{1}{1 - \rho_c} + x \left(\frac{1}{\rho_c} - \frac{1}{1 - \rho_c} \right) \quad (3.14)$$

$$\Rightarrow \bar{\tau}_c \in \left(\frac{1}{1 - \rho_c}, \frac{1}{\rho_c} \right). \quad (3.15)$$

From the definition of the density ρ , it is known that

$$0 < \rho_c < 0.5 \quad (3.16)$$

$$\Rightarrow 1 < \frac{1}{1 - \rho_c} < 2 \text{ and } 2 < \frac{1}{\rho_c} < \infty \quad (3.17)$$

and, depending on x , $\bar{\tau}_c$ can always assume values from an interval around the travel time of the homogeneous system, meaning that depending on x for any ρ_c the travel time can be smaller or larger than the travel time of the homogeneous system.

These results for the stationary TASEP with one bottleneck can be numerically verified for systems of lengths $L > 100$ and bottlenecks that are not too close to the initiation and termination regions to avoid edge effects. For the accuracy that I need for my statements later on, the distance has to be ≈ 10 sites away from the boundaries.

¹The interval for x is given by

$$x \in \left(0, \frac{1}{2} \left(1 - \sqrt{\frac{\bar{\tau}}{2 - \bar{\tau}}} \right) \right].$$

The interval for ρ is given by

$$\rho \in \left(0, 1 - \frac{1}{\bar{\tau}} \right].$$

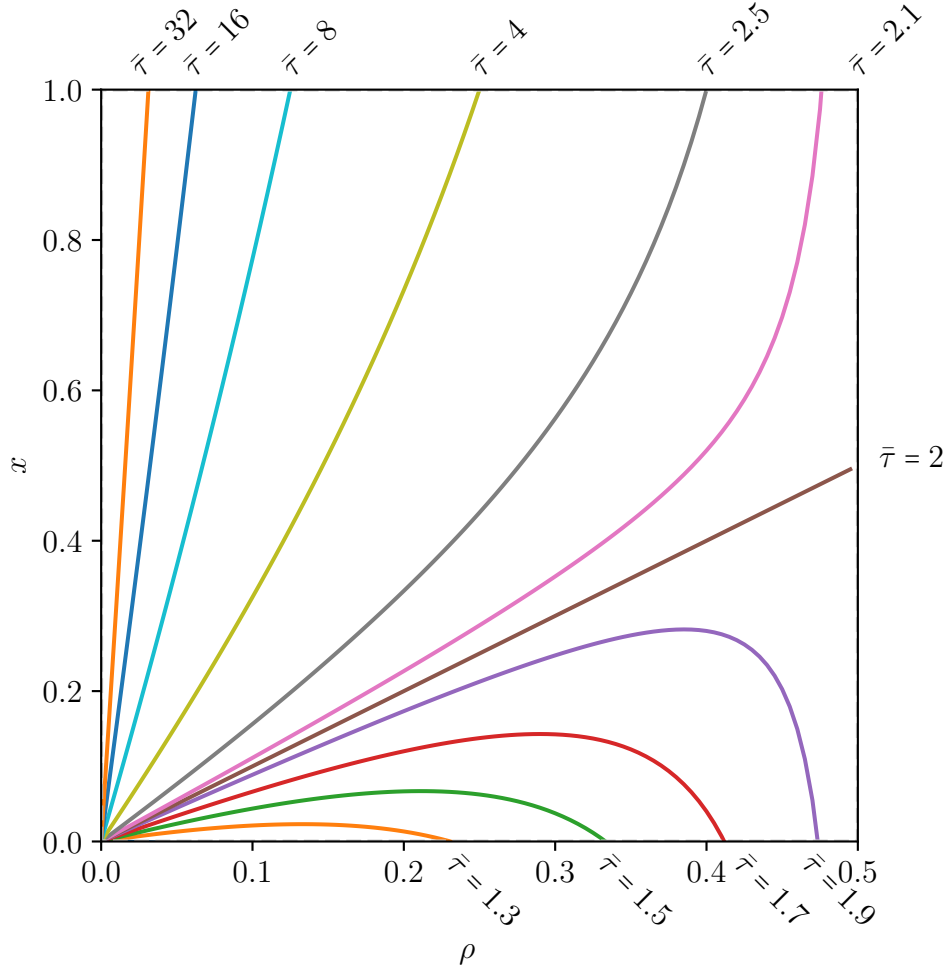


Figure 3.7.: This figure shows the phase space of the travel time $\bar{\tau}$. Each line represents a different value for the travel time $\bar{\tau}$. The homogeneous case is the line $\rho = 0.5$. The brown line, where $\bar{\tau} = 2$ represents systems, where the lower current in the system is exactly compensated by the lower average density. All lines start from $(0,0)$, but all but the line that compensates the current with the density approach the line $\rho = 0.5$, but never reach it.

The graphs were generated for different values of the travel time $\bar{\tau}$ by solving equation (3.13) for the relative location of the bottleneck x .

3.6.3. Interacting bottlenecks

When there are multiple bottlenecks present that are sufficiently far away from the boundaries and another, the strongest bottleneck, which is the one with the lowest rate, dominates the current J and density $\bar{\rho}$. The current J and density after the strongest bottleneck ρ are only dependent on the rate of this bottleneck, but not on the location. The average

3.6 Bottlenecks

density of the whole system $\bar{\rho}$ is not only dependent on the rate, but also on the location x of the strongest bottleneck, cf. equation (3.13). This is visualized in figures 3.8 and 3.9.

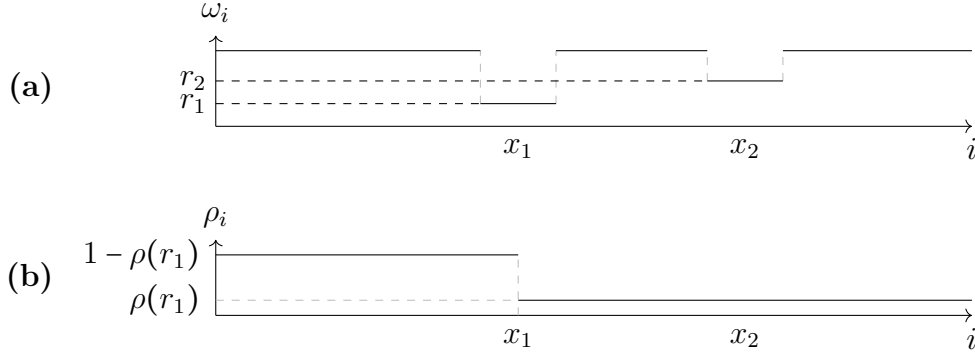


Figure 3.8.: An example system with two bottlenecks, r_1 and r_2 . r_1 is at position x_1 and r_2 at x_2 with $x_1 < x_2$ and $r_1 < r_2$. **(a):** The two bottlenecks in the system. $r_1 < r_2$, therefore the density ρ in **(b)** is only dependent on r_1 . The second bottleneck r_2 has no effects on the average density $\bar{\rho}$ or the current J .

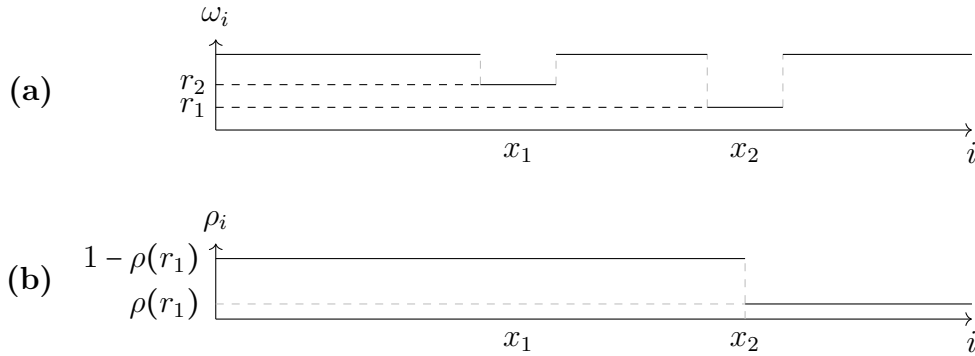


Figure 3.9.: The example from figure 3.8, but with swapped rates. Because $x_1 < x_2$ and $r_1 > r_2$ in **(a)**, the current is the same as before. But the average density changed and is now larger in **(b)**. The first bottleneck r_1 has no effects on the average density $\bar{\rho}$ or the current J .

Interacting bottlenecks are simple systems that can display sign epistasis (as described in section 2.4.1). This is easily visible when comparing the travel times $\bar{\tau}$ for all setups with two interacting bottlenecks in figure 3.6. In this example, the homogeneous wildtype F_0 has the density $\bar{\rho} = 0.5$, the current $J = 0.25$ and the travel time per site $\bar{\tau} = 2$. Two synonymous mutations change the respective rates to r_1 and r_2 at positions x_1 and x_2 span a landscape of 4 points with values of the measures from section 3.6.1.

I explain in a later section 4.1, that there is no sign epistasis in the current J , because then the whole landscape is only dependent on the rates r_1, r_2 , which are monotonic. If the measure chosen is the travel time, which is the average density divided by the current, the current is still monotonic in the rates, but the average density is not and therefore the travel time $\bar{\tau}$ is non-monotonic (cf. figure 3.6). Moreover, the travel time of interacting bottlenecks also shows neutral behavior. In figures 3.8 and 3.9, the larger bottleneck rate does not change the travel time of the system.

This interaction between the parameters is very interesting and is crucial for the model of interacting bottlenecks in chapter 5.

3.6.4. Relation between average density and bottleneck rate

There is no analytic function known for the relation between the average density $\bar{\rho}$ and the bottleneck rate in the system r . But the relation between the two can be calculated numerically. The result of this comparison is figure 3.10. One can calculate the density $\bar{\rho}$ of the TASEP with one bottleneck of rate r and find the value for the density ρ , because ρ is monotonic in r . This result can then be a map of densities to rates or rates to densities. This monotonic behavior is the reason why the rate and the bottleneck density can be used interchangeably.

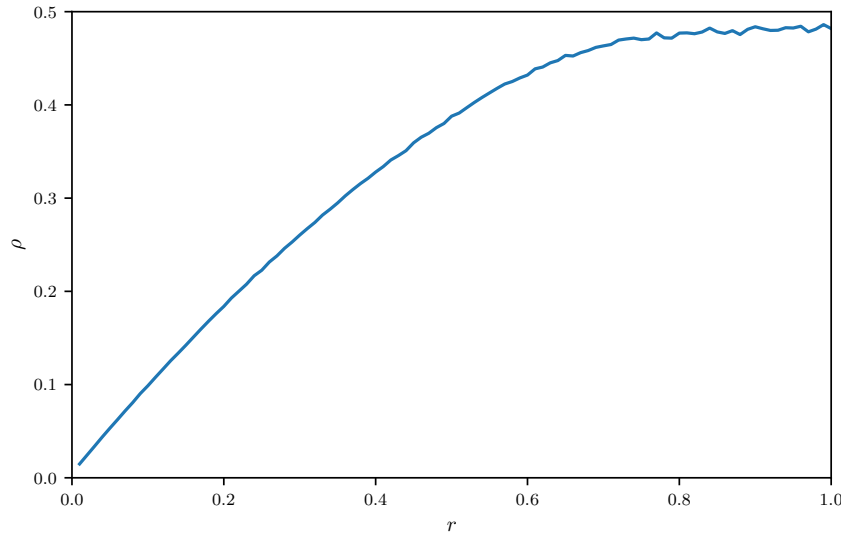


Figure 3.10.: The relation between ρ and r is numerically approximated. The statistical noise at the end comes from the close proximity to the maximum current phase.

4. (Non-)Monotonic parameters in a TASEP with bottlenecks

In the following, results from literature and numerical results are compared. The goal is to find a parameter, that has the required attributes of the fitness landscape. As it is the most commonly used parameter to describe fitness in I start with the current as a potential candidate for a fitness measure. There is an example from literature that supposedly shows sign epistasis in the current from the paper by Fouladvand et al. [7]. I show the part of their results that supposedly shows sign epistasis in section 4.1. Afterwards I explain a new and unpublished analytic proof by Krug [17] that disagrees with the example from the literature. To support the statement from the proof I analyze the more general case of random systems in section 4.2 both with an analytic approximation by Szavits-Nossan et al. [28] and with a numerical simulation.

The last section 4.4 shows that the phase space of two interacting bottlenecks features the intersections between regions where the mutations exhibit sign epistatic interaction and regions where they do not change signs, which is postulated in section 4.3.

4.1. The current is monotonic in the jump rates

In their paper Fouladvand et al. describe a TASEP with rates drawn from a binary distribution [7]. This method constructs a system with a given percentile of the system being small rates and the other part equal to 1. Two figures from the paper show a current that increases, if slow rates are added into the system (cf. figures 4.1a, 4.2a). The system described in the paper has a fast initiation rate $\alpha = 0.8$ and slow termination rate $\beta = 0.05$. The other rates in the system are fast $\omega = 1$ with probability $1 - f$ or have the rate $\omega = p_1$ with probability f . In figure 4.1a one can see an increase in the current, when comparing the homogeneous system at $f = 0$ (so the system where all rates are equal to 1) to the inhomogeneous system at $f = 0.1$ for all curves except for $p_1 = 0.05$. The other homogeneous system at $f = 1$ (here all rates are equal to p_1) always has a smaller current than at $f = 0.9$. This would directly proof, that there is sign epistasis in the current, because increasing the amount of smaller rates has an opposite effect if the background is different. The opposite slopes are visible around $f = 0$ and $f = 1$. This is not in accordance with the intuition of a TASEP systems current, which should be monotonic in the amount of slow rates.

When trying to reproduce these results, I find exclusively monotonic behavior (cf. figures 4.1b, 4.2b). The effect of changes on the particle current is always monotonic, meaning that an increase (decrease) of the jump rate ω_i at any site i in the system, causes the current to

4.1 The current is monotonic in the jump rates

increase (decrease) or stay the same. I want to stress that this example is just one of the results presented in the paper and it is also described as "unexpected" in there and that the extreme cases ($f = 0$, $f = 1$) in both mine and their graphs are in accordance with the homogeneous TASEP for the rates $f = 1$ or $f = p_1$ respectively.

If one adapts part II A. from Krug [16] to the TASEP system, there can not be any sign epistasis in the current of the systems described by Fouladvand et al. [7]. Even though the proof is concerned with surface growth in condensed matter systems, the arguments can be applied here as well¹. Furthermore, this proof does not only hold for the quenched spatial disorder system from Fouladvand et al., but is also true for general TASEP systems.

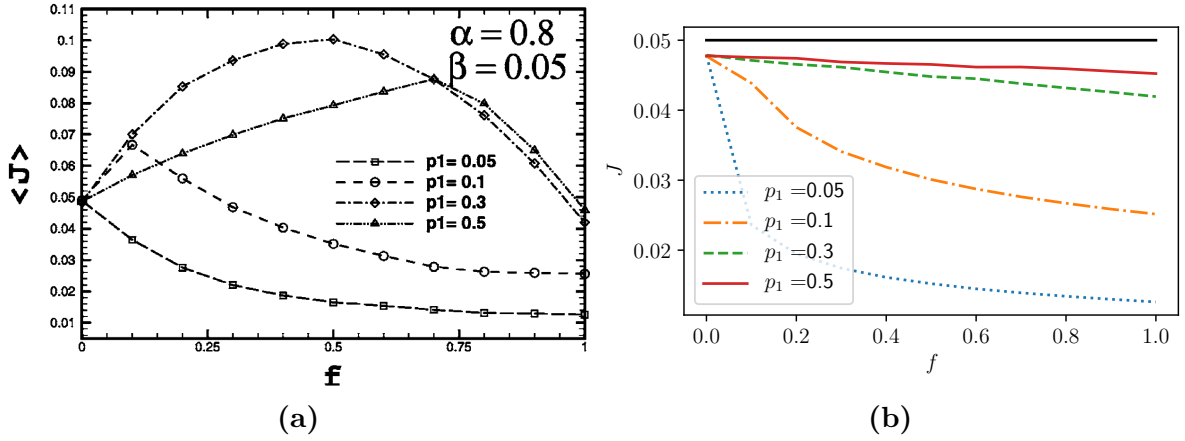


Figure 4.1.: Comparison between figure 8 from Fouladvand et al. [7] (a) and the same calculations with my own code (b). p_1 is the slower of two jump rates in the system. f is the percentile of sites that have rate p_1 . J or $\langle J \rangle$ is the average current of the TASEP. Both homogeneous cases, $f = 0$ (all rates are 1) and $f = 1$ (all rates are p_1) have the same values for J in both graphs but the behavior in the middle differs. In (a), the current is non-monotonic, in (b) the current is monotonic.

¹This part is based on unpublished notes by Krug [17].

4.2 Small fully random systems analyzed with a power series approximation and numerical simulations

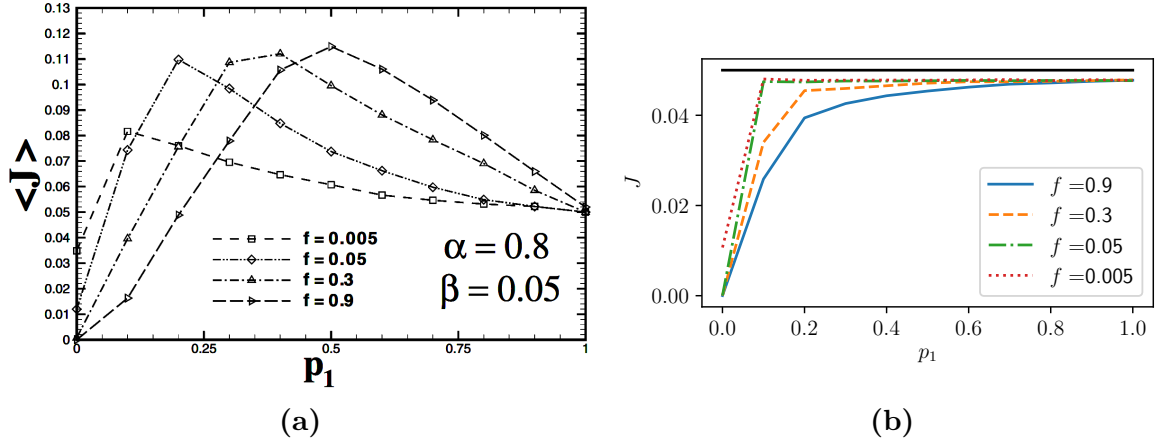


Figure 4.2.: Comparison between figure 14 from Fouladvand et al. [7] (a) and the same calculations with my own code (b). Here the parameter on the x-axis is the slower jump rate p_1 . Like in figure 4.1, the curves from the paper show non-monotonic behavior while the curves simulated by my code do not.

With this result from my work and with the proof by Krug, I conclude that there can not exist any TASEP current that is non-monotonic in its rates.

4.2. Small fully random systems analyzed with a power series approximation and numerical simulations

For systems with rates drawn from a distribution of random values, the current is a very complex function of the rates of the inhomogeneous TASEP. There is no general solution for it, but Szavits-Nossan et al. provide an analytical formula to compare to simulated data with their power series solution for the inhomogeneous TASEP with a small initiation rate α [28]. This is a result for a fully random system, with the constraint, that the initiation rate α is one order of magnitude smaller than the rest of the rates.

This is the analytic result, described by Szavitz-Nossan et al. [28] as the main result of their paper, for the expansion of the current J :

$$J(\alpha) = \alpha - \frac{1}{\omega_1} \alpha^2 + \left(\frac{1}{\omega_1} - \frac{1}{\omega_2} \right) \left(\frac{1}{\omega_2} + \sum_{j=3}^L \left(\frac{1}{\omega_j} + \delta_{j,L} \frac{1}{\omega_L} \right) \prod_{q=3}^j \frac{\omega_q}{\omega_1 + \omega_q} \right) \alpha^3 + \mathcal{O}(\alpha^4). \quad (4.1)$$

The epistasis measure described in equation (2.3) is used in the following analysis using a Mathematica code applying the equation (4.1). A small system size ($L = 4$) is filled with rates $\omega_i \in (0.1, 1)$. The initiation rate is chosen from the interval $\alpha \in (0, 0.1)$. The

4.2 Small fully random systems analyzed with a power series approximation and numerical simulations

two epistasis measures, E_1, E_2 , from equation (2.3) are dependent on four parameters (J_0, J_1, J_2, J_{12}) . There is sign epistasis if at least one of E_1, E_2 is smaller than zero. There is no sign epistasis, if both are larger than zero. Plotting the two epistasis measures against each other gives an overview of how many of the generated systems of rates $(\alpha, \omega_1, \omega_2, \omega_3, \beta)$ display sign epistasis in figure 4.3. The figure shows many examples of a negative value for E_1 or E_2 as red points. If there is no sign epistasis, the points are blue.

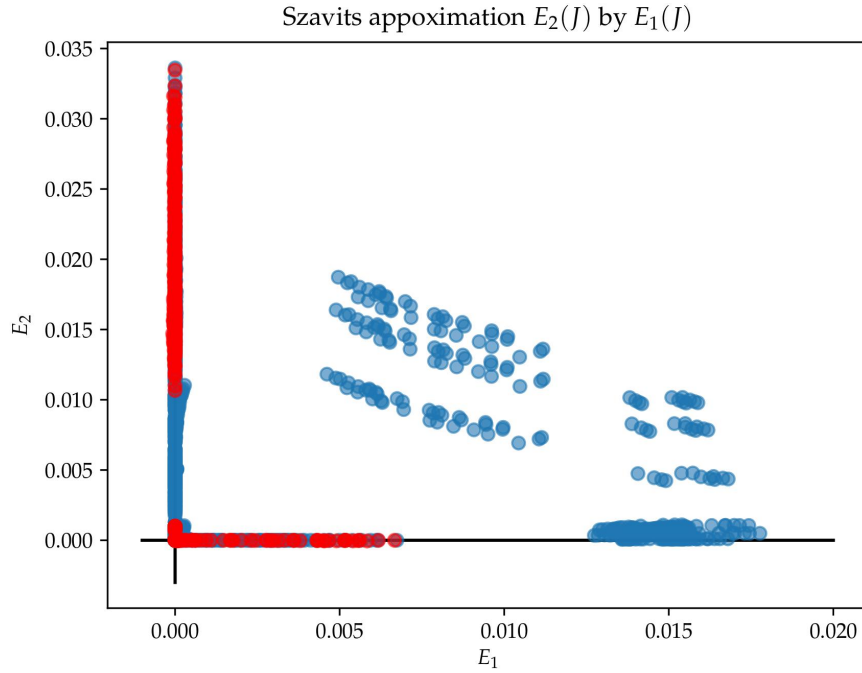


Figure 4.3.: E_2 by E_1 calculated with equation (4.1), system size $L = 4$, $\alpha = 0.5$, $\omega_i \in [1, 10]$. Each calculated system has a specific point in this landscape. If both E_2 and E_1 are positive the points are colored blue, if either of them are negative, the points are red. Many of the points close to the x-axis and y-axis are red, so the result is inconclusive whether or not sign epistasis exists in this case.

In the approximation, there are red points visible close to the lines $E_1 = 0$ and $E_2 = 0$. E_1 and E_2 is obtained for a large number of landscapes. For the red points, the minimal and maximal values are -10^{-4} and -10^{-11} respectively. The error due to the approximation of the parameters E_1 and E_2 is estimated to be $\mathcal{O}(\alpha^4)$, because the current itself is known up to order $\mathcal{O}(\alpha^4)$. For the values chosen, the error due to the approximation is of the same order as the negative values. This analysis does not show sign epistasis in the current, which is why for further analysis the results from the simulation are compared to

4.2 Small fully random systems analyzed with a power series approximation and numerical simulations

the numerical simulation results in the following section 4.2.1.

4.2.1. Calculating sign epistasis numerically

The results of section 4.2 are compared to results from numerical simulations. This section verifies results from section 4.3, gives motivation for the model described in chapter 5 and observes implications of it for real landscapes that are the topic of the next chapter.

By brute force calculation, the current J for systems of size $L = 4$ is calculated by calculating 10^5 landscapes which have 10 options for each rate α and ω at each site. E_1 and E_2 are calculated by choosing all combination of these landscapes to find possible sign epistatic effects of interacting random rates. All combinations of interacting systems are simulated with the numerical simulation explained in the appendix A. The results can be seen in figure 4.4. All points (E_1, E_2) are in the upper right quadrant, meaning that both $E_1 \geq 0$ and $E_2 \geq 0$. The error of the numerical calculation can be estimated as 10^{-9} to 10^{-11} by recording the current values and calculating the standard deviation of them over the time of simulation. But here the values are all positive.

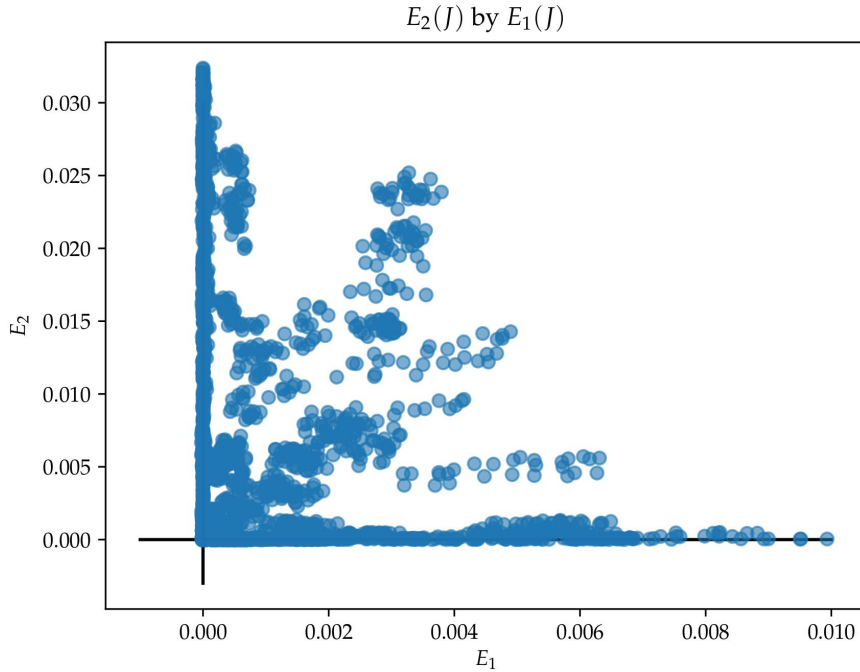


Figure 4.4.: E_2 by E_1 , System size $L = 4$, $\alpha = 0.5$, $\omega_i \in [1, 10]$. There are only points with positive E_1 and E_2 . This analysis does not show a single sign epistatic interaction.

4.3 Interfaces in the phase space where travel times are equal

This concludes the analysis of random landscapes with the insight that there is no evidence for sign epistasis in the current J of random landscapes. The results using the analytical result by Szavits-Nossan et al. [28] do not show sign epistasis beyond a reasonable doubt and the result from the numerical simulations do not show sign epistasis at all.

4.3. Interfaces in the phase space where travel times are equal

An interesting result of the model for interacting bottlenecks is, that it predicts the positions in the phase space of configurations for two interacting bottlenecks, where two travel times are equal. These are the points at which removing one bottleneck and adding the other does not have a fitness effect, or describes the brown line in figure 3.7, where adding a bottleneck to a homogeneous system keeps the travel time the same.

This phase space has dimension 4, because both the rates r_1, r_2 and the locations x_1, x_2 are a dimension in this. According to equation (3.13), the travel time $\bar{\tau}_i$ always depends on bottleneck position x_i and density ρ_i . One can compare two different systems by their travel times and therefore find the points at which the travel times are equal. This is compared to numerical simulations in section 4.4 and the lines in the phase space that are calculated here can be found in figure 4.5 and in the appendix B.

4.3.1. Interface where both bottlenecks have the same effect on travel time

$$\bar{\tau}_1 = \bar{\tau}_2$$

For two interacting bottlenecks at x_1 and x_2 with associated densities ρ_1 and ρ_2 , the travel times are equal, if

$$\bar{\tau}_1 = \bar{\tau}_2 \tag{4.2}$$

$$\Leftrightarrow \frac{1}{1-\rho_1} + x_1 \left(\frac{1}{\rho_1} - \frac{1}{1-\rho_1} \right) = \frac{1}{1-\rho_2} + x_2 \left(\frac{1}{\rho_2} - \frac{1}{1-\rho_2} \right) \tag{4.3}$$

$$\Leftrightarrow x_1 = \frac{\frac{1}{1-\rho_2} - \frac{1}{1-\rho_1}}{\frac{1}{\rho_1} - \frac{1}{1-\rho_1}} + x_2 \frac{\frac{1}{\rho_2} - \frac{1}{1-\rho_2}}{\frac{1}{\rho_1} - \frac{1}{1-\rho_1}} \tag{4.4}$$

for any fixed ρ_1 and ρ_2 , x_1 is a linear function of x_2 .

4.4 Comparisons of phase space interfaces to numerical data

4.3.2. Interfaces where a bottleneck has the same travel time as the homogeneous case

For equation (4.4) there is a solution $x \in (0, 1)$ for each value of $\rho \in (0, 0.5)$, where $\bar{\tau} = 2$. At these points in the phase space, adding the bottleneck to a homogeneous system does not change the travel time. It is, so-to-speak the line of neutral elements in the fitness landscape.

$$\bar{\tau} = 2 \quad (4.5)$$

$$\Leftrightarrow \frac{1}{1-\rho} + x \left(\frac{1}{\rho} - \frac{1}{1-\rho} \right) = 2 \quad (4.6)$$

$$\Leftrightarrow x = \frac{2 - \frac{1}{1-\rho}}{\frac{1}{\rho} - \frac{1}{1-\rho}} \quad (4.7)$$

$$\Leftrightarrow x = \frac{1-2\rho}{1-\rho} \cdot \frac{\rho(1-\rho)}{1-2\rho} \quad (4.8)$$

$$\Leftrightarrow x = \rho \quad (4.9)$$

The solution (4.9) is the point at which the smaller current, due to the insertion of a jump rate r smaller than the homogeneous rate, is compensated by the smaller average density. If $\rho < x$, the travel time $\bar{\tau}$ is larger with the slow bottleneck rate r than without it. If $\rho > x$, the particles in the system with the bottleneck travel faster. In the phase space there is one solution of this kind for $\bar{\tau}_1 = 2$ and for $\bar{\tau}_2 = 2$, the interface is a horizontal line for the rate r_2 inserted at x_2 and a vertical line for the rate r_1 inserted at x_1 .

When the travel times $\bar{\tau}_0, \bar{\tau}_1, \bar{\tau}_2$ are ordered by size, the interfaces in the phase space separate the regions, where the one travel times becomes larger than another. The comparison to numerical data is performed in the next section 4.4.

4.4. Comparisons of phase space interfaces to numerical data

The results from section 4.3 are compared to numerical simulations to find the predicted interfaces in the phase space.

For two interacting bottlenecks, r_1 at position x_1 and r_2 at position x_2 , the phase space has four free parameters, both locations x_1, x_2 and densities ρ_1, ρ_2 . It is useful to test if the intersections, described in section 4.3 are predicted correctly.

The phase space is scanned for configurations of travel times matching the conditions

4.4 Comparisons of phase space interfaces to numerical data

described in section 4.3. For each point $(x_1, x_2, \rho_1, \rho_2)$ in the phase space, four setups of the TASEP of length $L = 800$

- homogeneous,
- with jump rate r_1 at position x_1 ,
- with jump rate r_2 at position x_2 and
- with both rates present

are taken into account. Of these setups, the homogeneous system only has to be calculated once, and the systems with jump rate r_1 at position x_1 and jump rate r_2 at position x_2 can also be reused if already calculated. The densities $(\bar{\rho}_0, \bar{\rho}_1, \bar{\rho}_2, \bar{\rho}_{12})$ and the currents (J_0, J_1, J_2, J_{12}) are recorded. From these parameters, the travel times $(\bar{\tau}_0, \bar{\tau}_1, \bar{\tau}_2, \bar{\tau}_{12})$ are calculated. To avoid the need to calculate all points in the phase space, some combinations of ρ_1 and ρ_2 are fixed to specific values. I choose $\rho_1 < \rho_2$ and test, if $\bar{\tau}_{12} = \bar{\tau}_1$ as predicted in section 3.6.3. The results show, that for this large system ($L = 800$) there were no differences between the two values for the travel time that were larger than the error from transforming ρ into r as described in section 3.6.4.

The algorithm calculates new points in the landscape which are at distance $\Delta x = 0.025$ from the previous ones. Which points to calculate next depends on the results from the already calculated landscapes. Scanning the phase space like this, the algorithm finds the interfaces where two of the three travel times are equal. Notably the lines $\rho_1 = x_1$ and $\bar{\tau}_1 = \bar{\tau}_2$ are verified in all figures 4.5 and B.1-B.5.

The different colors mark how the travel times $\bar{\tau}_0, \bar{\tau}_1, \bar{\tau}_2$ are ordered. The lines, described in section 4.3.2, that are an analytic result of the section 4.3 are verified since the predicted gray lines match the intersections between the different colored stars.

This result is picked up in section 5.1 to give an argument that there is sign epistasis in some of the regions of the phase space, while others can not exhibit it.

4.4 Comparisons of phase space interfaces to numerical data

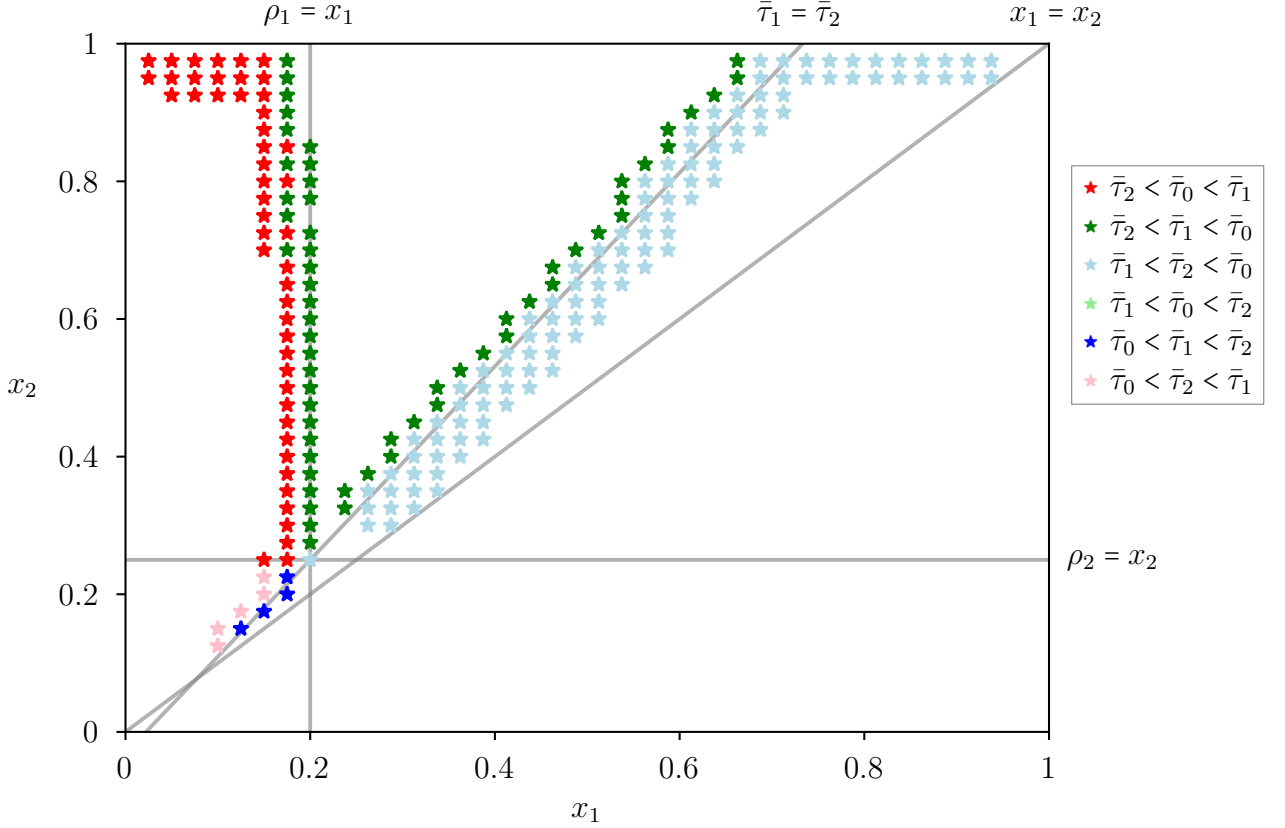


Figure 4.5.: The figure displays the phase space of x_1 and x_2 for parameters $\rho_1 = 0.2$ and $\rho_2 = 0.25$ with system size $L = 800$. One can see that the interfaces between points of different color fit to the theoretical predictions from section 4.3 (gray lines). The search algorithm starts at both the top right (lightblue stars) and the top left (red stars) of the phase space and then changes its direction when it encounters a change in the ordering of the three travel times (green stars). After the algorithm runs for too long or if it crosses the line $x_1 = x_2$ it stops. Note, that below the $x_1 = x_2$ line there was no calculation, because the points here are the same as the system above this line except that they have swapped ρ_1, ρ_2 values.

5. Modelling interacting bottlenecks

The last chapter showed that the current can not be the right fitness measure for a landscape with neutral mutations and sign epistasis. In the following I describe a model that uses the properties of the travel time as a fitness measure for a model of interacting bottlenecks, that takes both neutral mutations and sign epistasis into account.

The fitness landscape of antibiotic resistance is postulated to be a landscape of translation efficiency. This efficiency is the inverted travel time of the ribosomes on the *mRNA*. The shorter the travel time is, the faster the life-saving protein is produced for the cell and in the following the travel time is viewed as a fitness measure. For the description I use the expression fitness, but the assumptions come from the observations on the travel times in the TASEP from chapter 4.

The fitness is determined by the positions and rates of a system of interacting bottlenecks. For any combination of bottlenecks in this landscape, the bottleneck with the smallest rate dominates the current J and the density after the bottleneck ρ of the system (cf. section 3.6.3). The system of D interacting bottlenecks (r_1, r_2, \dots, r_D) , is then described by only two parameters, the bottleneck with the lowest jump rate r_i and the location x_i of this bottleneck. In figure 5.1a, the hypercube with four interacting bottlenecks is shown. This fitness landscape has the same dimension as the fitness landscape from Zwart et al. [35] for later comparison in chapter 6. This is a general hypercube as used by Wright [33], that my model has not been applied to yet. The landscape is dependent on all 4 rates (r_1, r_2, r_3, r_4) and 4 locations (x_1, x_2, x_3, x_4) of the 4 bottlenecks. Because the bottlenecks do not have distinguishable features aside from their rates and location, without loss of generality they can be ordered for ease of notation as

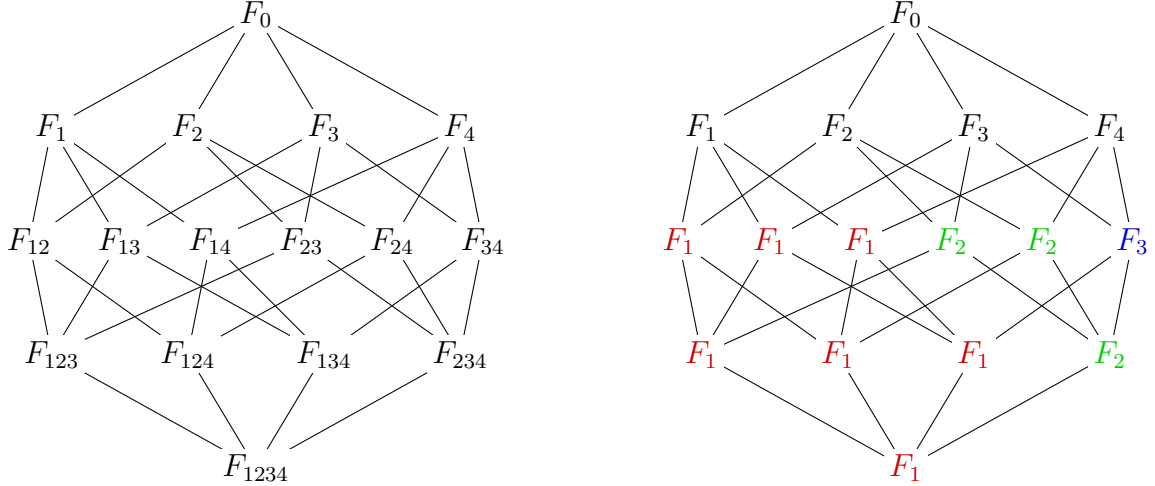
$$r_1 < r_2 < r_3 < r_4 < 1 . \quad (5.1)$$

It is important, that even if the rates are ordered, the locations x_1, x_2, x_3, x_4 can still be chosen freely. All systems with fixed rates can be displayed in this way, because with the right choice of locations, any sorting of fitness values can be created. This follows from equation (3.17) because for any fixed r_1, r_2, r_3, r_4 , there is a set of x_1, x_2, x_3, x_4 so that the fitness values F_1, F_2, F_3, F_4 and the value $2 = F_0$ can be in any order. The system is displayed as a hypercube of four dimensions in figure 5.1a.

Starting from figure 5.1a, because the strongest bottleneck dominates the travel time, whenever the bottleneck r_1 is in the system, the fitness is equal to F_1 . When r_1 is not present, but r_2 is, the fitness is equal to F_2 . Following this logic, the number of individual

CHAPTER 5. MODELLING INTERACTING BOTTLENECKS

fitness values are reduced from 2^D to $D + 1$ with D the dimension of the hypercube, corresponding to the number of bottlenecks. After applying this key assumption of my model, the system changes from figure 5.1a into figure 5.1b.



(a) The Wrightian fitness landscape of dimension 4 as in the original paper [33]. The points are associated with the fitnesses and the lines are the mutations.

(b) The fitness landscape from (a) after applying my model. The red colored fitness values are dominated by the bottleneck r_1 and therefore become F_1 . The green colored fitness values are dominated by the bottleneck r_2 , because the rate r_1 is not present and the blue colored fitness value is dominated by the bottleneck r_3 . The fitness of the homogeneous system (F_0) and the system with the weakest bottleneck (F_4) does not change under the assumption.

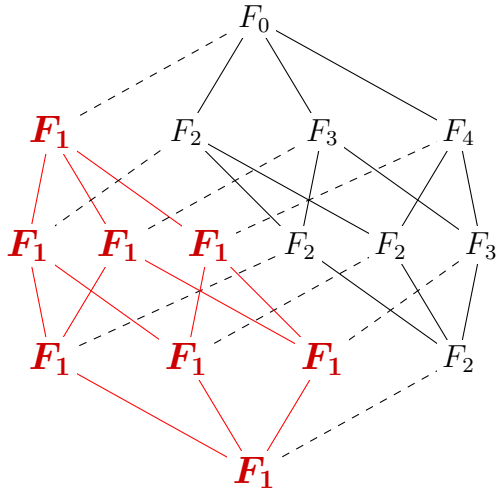
Figure 5.1.

In essence, because the rates are ordered as in (5.1), the dominant fitness effects are ordered as well. This method reduces the number of different fitness values from 16 to 5. In figure 5.1, there are "subcubes" within the hypercube. These are hypercubes of lower dimension contained in the landscape of higher dimension. A hypercube of dimension 4 can be "cut" into two hypercubes of dimension 3, displayed in figure 5.2a. One can see that the red cube in this figure only features the fitness values F_1 and the black cube does not contain F_1 at all. This is because the connection between the two subcubes is the mutation that adds (or removes) r_1 into (or from) the system. If r_1 is present, the fitness value is equal to F_1 . Within the red cube the fitness landscape is flat, meaning that the fitness values are the

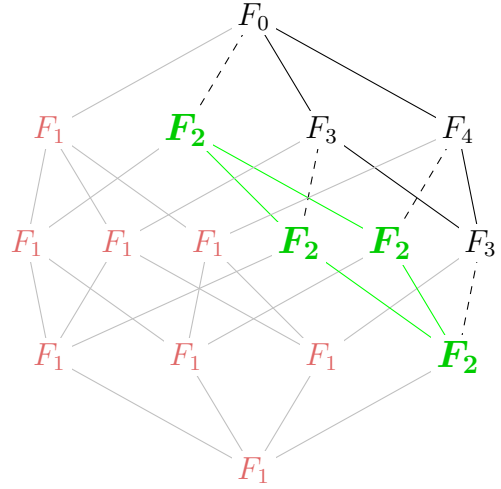
same, even though the microscopic configurations of jump rates are different at every point in the landscape. The black cube in figure 5.2a has the same shape as if the mutation r_1 was not part of the landscape in the first place and one can perform a second cut, displayed in figure 5.2b, similar to the one in figure 5.2a.

This is the description of my model of interacting bottlenecks to explain fitness landscapes where the travel time is the fitness measure. It features many neutral (or flat) mutations and sign epistasis due to dominant bottlenecks, that interact with each other.

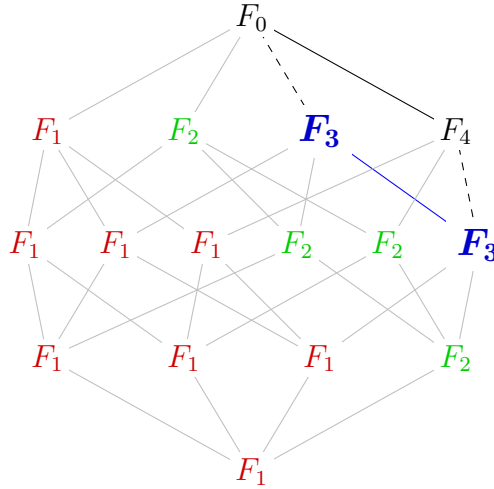
At any point of the cutting process, one of the subcubes has to have only flat mutations, meaning that all mutations within this subcube have no fitness effect. I call this subcube a flat subcube, because all mutations within it are flat as explained in section 2.3.1. This is the main assumption of my described model and is used as a testing mechanism in the next chapter 6 to find out if my model accurately describes the experimental data. All orderings of fitness values can be obtained if one is free to choose the rates here, as described in section 3.6.2.



(a) When the hypercube is cut along the dashed lines, it splits into two hypercubes of three dimensions. Note that the fitness values on the red cube are all F_1 and on the black cube none of them is F_1 .



(b) The black subcube of dimension 3 from (a) can be cut again to form two subcubes of dimension two. The cut is performed along the dashed lines. Again, the half of the landscape that is dominated by the strongest bottleneck, which is r_2 here, is cut off.



(c) The black subcube of dimension 2 from (b) can be cut once more. This leaves the black and blue lines in this figure. The black nodes F_0 and F_4 , do not matter for my model, because they may take any value. Because these two values are not restricted by my model, they do not change the applicability of it. The original system of 2^D fitness values is now described by only five fitness values F_0 (black), F_1 (red), F_2 (green), F_3 (blue) and F_4 (black), that cover all nodes in the system.

Figure 5.2.

5.1. Types of two-dimensional subcubes

I explain a feature of the two dimensional subcubes within the landscape from figure 5.2c. I call the two rates that span the subcube r_1 and r_2 and therefore the fitness values are F_0 , F_1 , F_2 and F_{12} . By choosing $r_1 < r_2$ it follows that $F_{12} = F_1$. If there is a background on which the mutations occur, then there are three types of interactions between the two mutations on the subcube.

If there exists a rate r_{dominant} in the subcube, that is smaller than r_1 , so $r_{\text{dominant}} < r_1 < r_2$, then the landscape is fully neutral and all fitness values at the corners are the same.

If there exists a rate $r_{\text{semi-dominant}}$ in the subcube, that is smaller than r_2 , but larger than r_1 , so $r_1 < r_{\text{semi-dominant}} < r_2$, then the fitness landscape looks like in figure 2.3a, the weaker of the two mutations does not have a fitness effect, so it is neutral, but the other one does have a fitness effect.

If there two rates are the smallest two in the system, then the fitness landscape can exhibit sign epistasis. The reason for this can be found in figure 5.3. It shows that there are two ways of sorting the fitness values F_0, F_1, F_2 in the two dimensional fitness landscapes that always exhibit sign epistasis and the others do not. This result gives boundaries for the possible interactions in the 2-dimensional subcubes which feature a set of three different fitness values. In this case, one of the orderings of travel times in 5.3 is present.

The fitness values for given values of r_1, r_2 can change depending on the positions at which they are inserted into the system (explained in section 3.6.2). The six landscapes displayed in figure 5.3 are all possible orderings of F_0, F_1 and F_2 for a background on which r_1 and r_2 are the bottlenecks with the lowest rates.

The figures 4.5 and B.1-B.5 connect the locations of the bottlenecks to the existence of sign epistasis, given that the densities are known. One could also calculate the phase space of the two densities given that the locations are known this way.

5.1 Types of two-dimensional subcubes

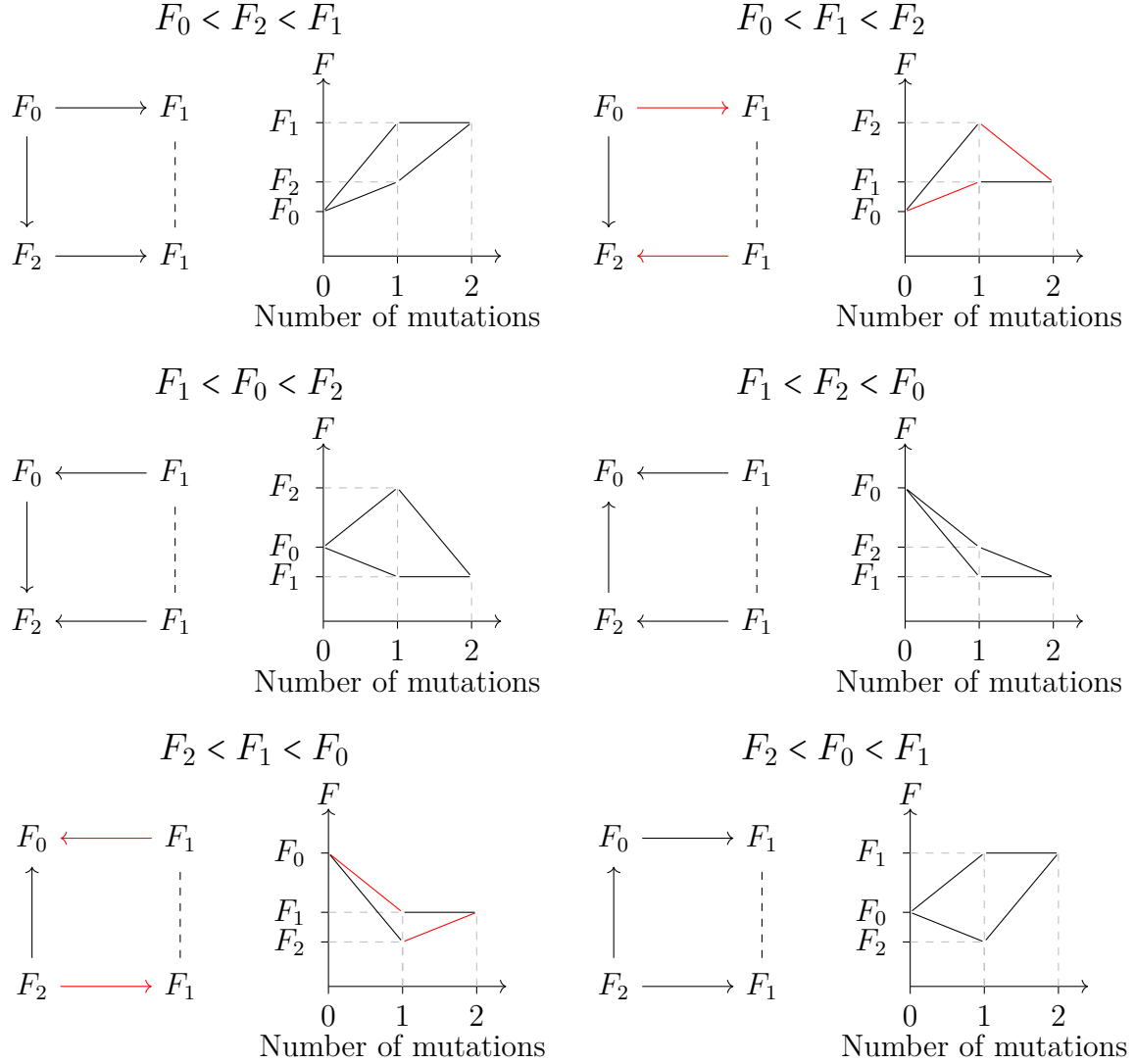


Figure 5.3.: F_0 , F_1 and F_2 can be ordered in a total of six different ways. This leads to six possible interactions between the different fitness values. Displayed as a hypercube (middle), the arrows point to larger fitness values (and therefore smaller fitness). Recalling section 2.4.1, arrows pointing in opposite directions show sign epistasis (red arrows). In most experimental papers these landscapes are displayed like on the right. In section 6.1a I search for these 2 dimensional subcubes in the results from an experiment by Zwart et al. [35].

5.1.1. Inhomogeneous case

It is not necessary that F_0 is the fitness value of the homogeneous system. It can rather be the fitness values of the system, where the mutations m_1 and m_2 did not occur yet. Mutations can not only add slow rates into the system, but also remove them. It is imag-

5.1 Types of two-dimensional subcubes

inable, that the rate r_1 is already present in the system, and that mutation m_1 removes it from the system, by replacing it with a faster rate. Because F_1 is the fitness value if the rate r_1 is present and r_2 is the fitness value if r_2 is present, but r_1 is absent. If the rate r_2 is added to the system, it has no effect if $r_1 < r_2$. Another mutation m_1 removes the slowest bottleneck rate r_1 from the system, leaving the rate r_2 as the strongest bottleneck rate. Then the landscapes from figure 5.3 are "flipped". The flat part of the landscape is the connection between the wildtype F_1 and the landscape with both slow rates present, which also has the fitness value F_1 .

It is important to note, that these ways of interacting mutations are the only possibilities for interacting bottlenecks in this model. For all landscapes described by this model, one of the mutations in a 2-dimensional subcube needs to be flat, or else the model can not describe the interaction.

6. Analysis of the Zwart et al. landscape

In their paper, Zwart et al. analyze a set of experimental data for fitness effects of synonymous mutations [35]. The measure for fitness is the concentration of the antibiotic cefotaxima, at which 99.99% of cells die. This measure is defined as the $IC_{99.99}$, which is defined in a paper by Schenk et al. from 2012 [24]. The species that survives in higher concentrations of the antibiotic therefore have a higher $IC_{99.99}$ value. In the Schenk et al. paper [24], the fitness effects of 48 synonymous and non-synonymous mutations are already discussed and in their 2013 paper [25], a fitness landscape is constructed.

The Zwart et al. paper [35] examines the fitness landscape of synonymous mutations more thoroughly, shows the interactions between them and compares the single effects quantitatively. All of the synonymous mutations observed are constructed in a lab and the $IC_{99.99}$ is calculated with a so-called assay, which are media with different concentration of the antibiotic in which the cells grow. The cell colonies are counted after they have grown and the $IC_{99.99}$ is calculated. There are three replicates for each species to find the error of measurement of the $IC_{99.99}$.

The results of the Zwart et al. paper [35] show surprising results. Synonymous mutations affect the fitness of the cells and some synonymous mutations have a much larger effect on fitness than non-synonymous mutations. This is against the intuition, that a mutation, that does not affect the constituents of the protein, still has a fitness effect for the cell.

In this chapter I compare the model from chapter 5 to the experimental data and show an interesting feature of the fitness landscape that has not been described by Zwart et al. [35]. I start by explaining what the implications of the model from chapter 5 are for experimental data in section 6.1.

6.1. Do the assumptions apply to experimental data?

For the model from chapter 5 to accurately describe the data, the following assumption must be true. At any point of the cutting process, described in chapter 5, half on the system must be a flat subcube. Whether or not this assumption is applicable to experimental data of D mutations, can be tested with the 2^D fitness values of the different species within the mutation landscape.

If the experimental fitness values are known, the hypercube of dimension D is cut. There are $D^2 D!$ ways to cut the hypercube into a set of $D-1$ subcubes of dimension $D-1$, $D-2$, ..., 1.

6.2 Applying my model

The deviation from a flat landscape, where all fitness values are the same, ΔF is calculated by adding up the standard deviation of the elements of the set of subcubes

$$\Delta F = \sum_{m=0}^D \sigma(F_m) . \quad (6.1)$$

There are other ways of weighing the standard deviations within the subcubes possible. For example the deviation of the larger cubes could be weighted more than those of smaller cubes. In the following example, I choose to equally weigh all subcubes. This makes sense for testing if the subcube is flat, because I want to know that each subcube has a constant value and not necessarily that each fitness value has the same value as the average value of its respective subcube. This is a small detail though, because for the landscape from Zwart et al. [35], the best fitting cut is the same when calculated both ways.

I want to emphasize, that testing if the subcubes are flat, is done solely with the fitness values. The locations x_1, x_2, \dots, x_D or the rates r_1, r_2, \dots, r_D do not need to be known. Because the average fitness within the subcube is the characteristic feature of it, I use F_1, F_2, F_3 for the average fitness values of the subcube and for identifying the subcubes.

For a hypercube of 4 dimensions, there are 192 sets of subcubes (one of dimension 3, 2 and 1 each). After calculating the 192 sums of the standard deviations (6.1) for the three subcubes, I choose the cut with the smallest deviation ΔF in the next section 6.2 to investigate if my model describes this cut of the fitness landscape.

6.2. Applying my model

The experimental fitness landscape of Zwart et al. [35] is shown in figure 6.1a. There are both positive (black lines) and negative (red lines) fitness effects visible in the landscape. The error of measurement given in the paper is shown as bars at each fitness value. There are many mutations which do not have a strong fitness effect, shown as lines with a smaller slope than the error of measurement. These are assumed to be neutral, meaning that they do not have any fitness effect. Zwart et al. [35] classify the fitness values into three categories, which are small, medium and large fitness values. The categories however, do not capture the structure of the mutation landscape. If different species with the same fitness values are connected by neutral mutations needs to be checked.

In the following, a more detailed look is taken at the $IC_{99.99}$ fitness values with regard to sign epistasis, as described in section 2.4.1 and neutral mutations, as explained in section

6.2 Applying my model

2.3.1.

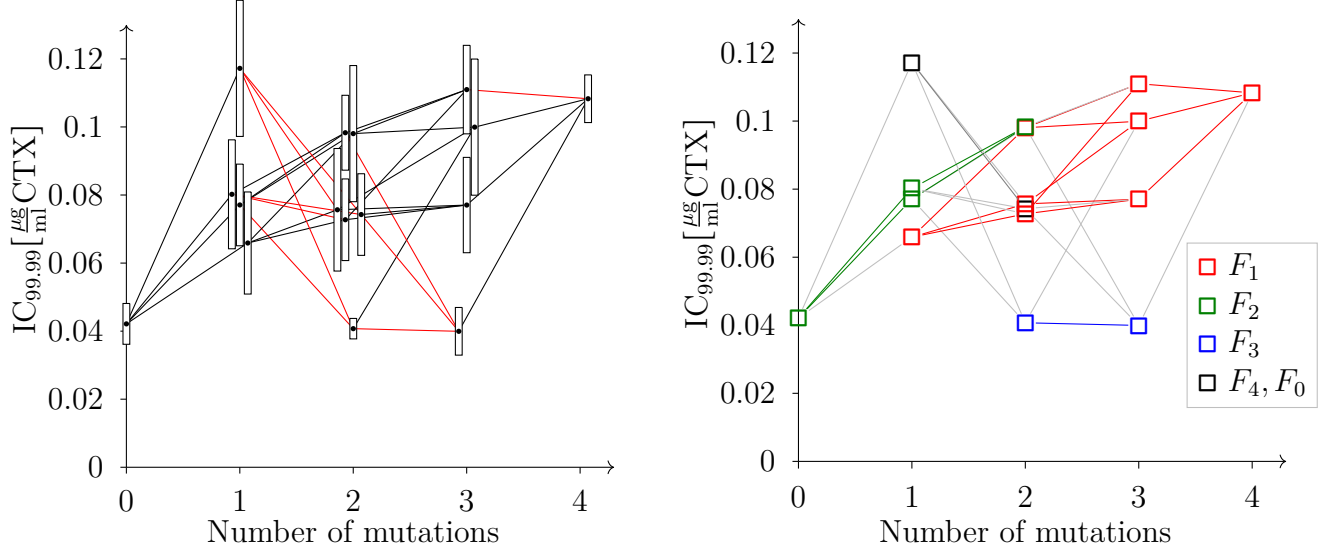
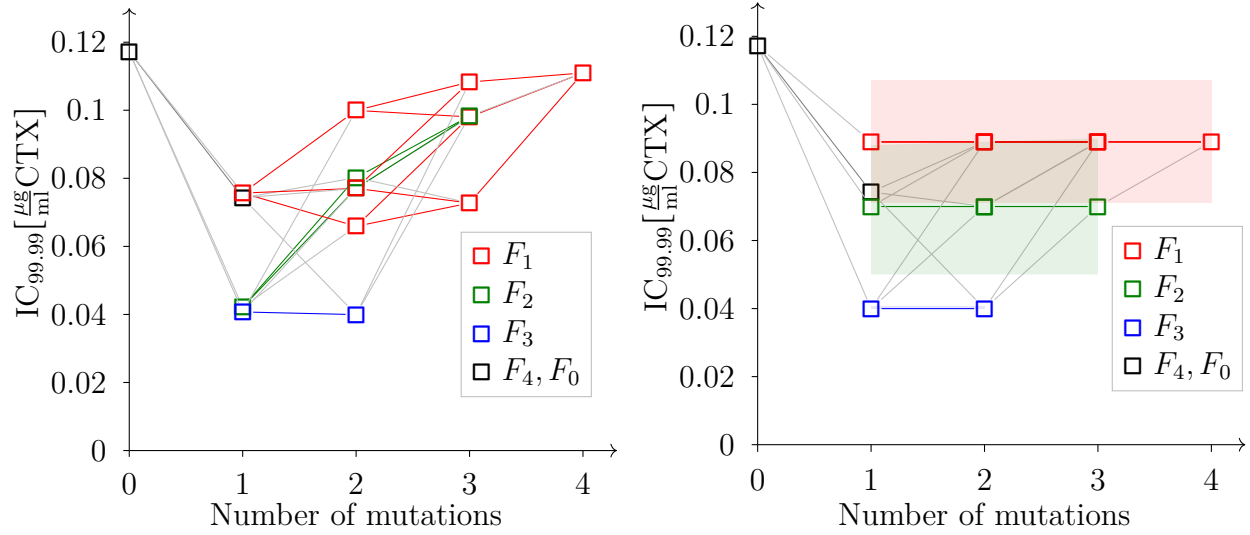


Figure 6.1.

For the Zwart et al. [35] landscape, I cut the best fitting subcubes as described in section 6.1. After the fitness values are categorized by the subcube they belong to, the landscape can be displayed as in figure 6.1b. One can see, that the fitness values within a subcube (points of the same color) are similar, but not the same. If the model from chapter 5 described the landscape perfectly, there must be one way of cutting the landscape into flat subcubes. The averaged experimental values for F_1, F_2 and F_3 with their standard deviations are shown in table 6.1 and visualized in figure 6.2b. The resorting in figure 6.2a is performed to mimic the structure of figure 5.2c, where an origin F_0 is the point from which the hypercube spans over the space of four mutations. This resorting leaves

the structure the same and is just a visual aid to help understanding the structure of the hypercube.



(a) The same data as in 6.1b, but resorted to better visualize the assumptions of the model. I choose the new origin of the hypercube at the species with only the mutation 87* present (cf. section 5.1.1). All other fitness values move one step closer to the origin if they contain the mutation 87* and one step further away if they do not contain this mutation.

(b) If the assumptions of chapter 5 apply perfectly, the landscape would look like this figure. The average value is shown as strong colors and the standard deviation within each subcube is shown as translucent rectangles around the values. I want to emphasize, that these rectangles do not represent the error of measurement.

Figure 6.2.

In figure 6.2, one can see that there are large differences between the values expected by the model and the experimental data. If the assumptions drawn by the model were correct, the subcubes in figure 6.2a of the same color should have the same fitness values. Then their standard deviation would be equal to zero in figure 6.2b.

The standard deviations $\sigma(F_m)$ are compared to the minimal experimental uncertainty $\Delta\text{IC}_{99.99}$ in the subcube. This is the most strict measure to see if the model fits to the experimental data. The results are shown in the following table.

6.3 Additive mutational effects within subcubes

F_m	Average $IC_{99.99}$	$\sigma(F_m)$	Minimum of $\Delta IC_{99.99}$	Is the subcube flat?
F_1	0.089	0.018	0.007	not flat
F_2	0.07	0.02	0.006	not flat
F_3	0.0404	0.0004	0.003	flat
Applying equation (6.1), the total deviation is $\Delta F \approx 0.04$				

Table 6.1.: Table of the average fitness values (Average $IC_{99.99}$) within the subcubes F_1, F_2, F_3 and their standard deviation $\sigma(F_m)$. If the standard deviation is smaller than the minimum of the error of measurement $\Delta IC_{99.99}$, the subcube is flat. This is only true for F_3 . The total deviation from my model is $\Delta F \approx 0.04$, which is much larger than the experimental errors of measurement.

Only the subcube F_3 fits to the model, so the model does not describe the landscape. The subcubes F_1 and F_2 are not flat. To better represent the experimental data, the model is adapted to account for the non-flat subcubes in the next section 6.3.

6.3. Additive mutational effects within subcubes

The last section showed that the mutational landscapes are not as simple as predicted. The result in figure 6.2a shows another feature of the subcubes. Within the subcubes F_1 and F_2 , the mutations seem to be additive, a feature of fitness landscapes described in section 2.4.2, which means that mutations have a constant fitness effect. Additive landscapes like this are characteristically described by the current J of the system and not the travel time $\bar{\tau}$. If the subcube landscape is additive, the synonymous mutation landscape is mainly dominated by the translation efficiency (which is the inverted travel time). Within the subcubes of equal travel time, the behavior can be different. The new assumption of this section is, that there may be some mutations, that increase the current within the subcubes of equal travel time and therefore the amount of protein produced, resulting in an increase of the fitness. This assumption combines the traditional measure for fitness, which is the current (because more antibiotic-digesting protein is better for the organism), with the newly proposed measure, which is the translation efficiency (because faster production of the antibiotic-digesting protein is beneficial for the cell). This new assumption is only based on the visible additive effects of certain mutations in the landscape in figure 6.2a. The biology behind it could be very different and needs to be discussed with researchers that have more expertise in the biological side of the process. Whether or not this new assumption is realistic is an open question and this alteration of the model should therefore be regarded as a conjecture rather than a fact.

6.4. Comparing my adapted model to the experimental fitness landscape

The previous section 6.3 gives arguments for linear effects inside of subcubes. This section analyzes the experimental data and shows that the experimental fitness landscape can be describes with some additions to the model.

In the following, I use the expression *mutational effect* E_n to describe the difference between the experimentally measured fitness values $F(a), F(b)$ of two species a, b that are connected by a mutation c in the fitness landscape. To be precise, mutation c on species a transforms it into c and the mutational effect on fitness is $E_c = F(b) - F(a)$.

The averages \bar{E}_n ¹ and standard deviations $\sigma(E_n)$ of the mutational effects, which are the absolute fitness differences between two species connected by a mutation n within a subcube F_m , are calculated. The standard deviation $\sigma(E_n)$ is compared to the minimum of the experimental errors of the mutational effects $\Delta_{\text{exp}}(E_n)$. The experimental errors $\Delta_{\text{exp}}(E_n)$ for the effects can be calculated for each mutation in the subcube as the errors of the fitness values along the edges. If a mutation c connects the points a and b , with the experimental errors Δa and Δb , then the error along the edge is $\Delta_{\text{exp}}(E_c) = \sqrt{\Delta a^2 + \Delta b^2}$. Because I choose the minimum of the experimental errors for each mutation, I get the lower bound for the mutation having a significant fitness effect. I underestimate the error of measurement on purpose to be sure to not classify a mutation which has an additive effect as flat. This gives a comparison of the variation of the fitness effects caused by the mutations and the equivalent error of measurement $\Delta_{\text{exp}}(E_n)$.

Table 6.2 shows, that the deviations $\sigma(E_n)$ are smaller than the errors of measurement $\Delta_{\text{exp}}(E_n)$. The mutation 9^* in the subcube F_1 and the mutations 9^* and 89^* in subcube F_2 are additive because their average effect $\|\bar{E}_n\|$ is very large compared to the standard deviation $\sigma(E_n)$, so they add a large constant value to the fitness. They are also significantly larger than the experimental error $\Delta_{\text{exp}}(E_n)$.

The mutations 87^* and 89^* in subcube F_1 have average mutation effects \bar{E}_n that are much smaller than the experimental error of measurement $\Delta_{\text{exp}}(E_n)$. Therefore the maximal mutational effect $\max(\|E_n\|)$ is calculated to compare it with the minimum experimental error $\Delta_{\text{exp}}(E_n)$. If the experimental error is smaller, this proofs that the mutation has no

¹This is the average of the mutational effect E_n over all fitness values that are connected by mutation n within the subcube.

6.4 Comparing my adapted model to the experimental fitness landscape

effect and needs to be considered flat.

The subcube F_3 is also shown in table 6.2 for completeness. In last section 6.2 this mutation is already shown to be flat.

Subcube F_1	$\ \bar{E}_n\ $	$\sigma(E_n)$	$\Delta_{\text{exp}}(E_n)$	$\max(\ E_n\)$	Effect in landscape
Mutation 9^*	0.031	0.005	0.014	–	additive
Mutation 87^*	0.003	0.004	0.013	0.010	flat
Mutation 89^*	0.007	0.004	0.019	0.013	flat
Subcube F_2	$\ \bar{E}_n\ $	$\sigma(E_n)$	$\Delta_{\text{exp}}(E_n)$	$\max(\ E_n\)$	Effect in landscape
Mutation 9^*	0.027	0.008	0.013	–	additive
Mutation 89^*	0.030	0.009	0.016	–	additive
Subcube F_3	$\ \bar{E}_n\ $	$\sigma(E_n)$	$\Delta_{\text{exp}}(E_n)$	$\max(\ E_n\)$	Effect in landscape
Mutation 89^*	0.0008	–	0.004	–	flat

Table 6.2.: The average absolute mutational effect $\|\bar{E}_n\|$ is compared to its standard deviation $\sigma(E_n)$ and the experimental error $\Delta_{\text{exp}}(E_n)$. This shows that mutations 9^* in subcube F_1 and 9^* and 89^* in subcube F_2 are additive. They have a significant effect on the fitness in their respective subcubes. The mutations 87^* and 89^* in subcube F_1 are flat, because they have a smaller maximal mutation effect $\max(\|E_n\|)$ than the experimental error $\Delta_{\text{exp}}(E_n)$. Their effect on fitness is regarded as 0. Mutation 89^* from subcube F_3 is added for completeness, but it was shown to be flat in the last section 6.2 already. Values that were not necessary to calculate are omitted from this table.

The results of this section are collected in table 6.2. It shows that all mutations within subcubes of the Zwart et al. [35] fitness landscape are either additive or flat. Allowing this additive substructure in my model therefore explains all fitness effects in the landscape. The fitness landscape is visualized in figure 6.3, which can be compared to figure 6.2a.

6.4 Comparing my adapted model to the experimental fitness landscape

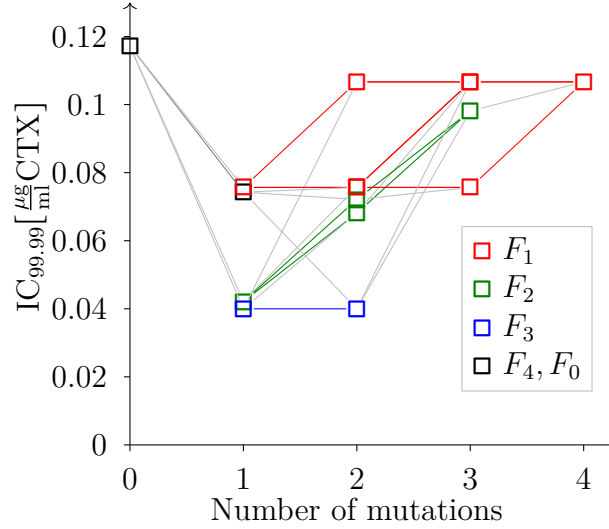


Figure 6.3.: The adaptation of the model improves the theoretical landscape significantly. Different to figure 6.2b, the adapted model from section 6.3 considers additive mutation effects, it can describe the experimental data much better. The structure of the fitness landscape is very similar to figure 6.2a and the model can fully replicate the experimental landscape.

Figure 6.3 recreates the original landscape 6.2a well. This could lead to a better understanding of the landscape of synonymous mutations. The large value for the deviation from my model from chapter 5 ΔF from equation (6.1) in section 6.2 can be attributed to additive fitness effects in the subcubes. The reason for this behavior could be that mutation 9^* is in the initiation region and therefore changes the current differently than the rates that are further away from the boundaries of the system. Furthermore, the mutations 87^* and 89^* might interact on a short range due to effects described in the section about edge effects 3.4.2. Bottlenecks in the initiation region and bottlenecks that are close together are not immediately described by my model from chapter 5.

The fitness is dominated by the travel time, but if bottlenecks are close together or close to the boundaries of the system, the current has an additive effect on the fitness. I conclude that the overall fitness in the landscape is a function of both the travel time and the current. It can as such can be described with my model from chapter 5 if the mutations are sufficiently far apart from each other and also far from the boundaries.

7. Conclusion and discussion of results

The experimental fitness landscape analyzed in this thesis exhibits sign epistatic effects and many neutral mutations. This is visible in figure 6.1a. The primary approach of using the current as a fitness measure does not yield correct results, because current of a TASEP system can not exhibit sign epistasis. This is shown in section 4.1 where I explain, that the current can not explain the sign epistatic effects in the fitness landscape of synonymous mutations with a recent proof by Krug combined with my analysis of TASEP landscapes with random rates.

Because the current is not the right measure for fitness, the travel time, as described by Szavitz-Nossan et al. [28], is analyzed for its properties. The homogeneous TASEP system with inhomogeneous rates, called bottlenecks, is analyzed for its properties in section 3.6. The analysis shows that a homogeneous TASEP with two bottlenecks can already exhibit sign epistasis in section 3.6.

From this observation, I construct a model that features multiple bottlenecks in a TASEP in chapter 5. This model has the travel time as the only measure for fitness and the prerequisite, that the current is monotonic in the bottleneck rates. The phase boundaries in a TASEP with interacting bottlenecks are shown to be well represented using results from my model in section 4.4 and are numerically verified in section 4.2.1. Finally, the model is applied to the experimental data from Zwart et al. [35], that initiated this search for a model of synonymous mutations in chapter 6.

The analysis of the experimental data shows, that the general structure of the experimental fitness landscape can be recreated correctly with my model, but it can not explain all effects. The subcubes of the landscape are predicted to be fully neutral, but display non-neutral behavior in section 6.2. This non-neutrality is thereafter shown to be a monotonous effect. This effect is likely the result of a change in the current within the subcubes, which is contributing as a second factor to the fitness of the organism as conjectured in section 6.3. This result is an extension of my model from chapter 5.

Idealized systems of interacting bottlenecks, which are far enough apart from each other and the boundaries, do not display edge effects and are therefore well described by my model from chapter 5. The experimental data that was analyzed is not one of these systems as it contains a bottleneck near the initiation region and two bottlenecks that are only one codon apart. Because the general description fails to describe the experimental data, the

model is adapted by taking into account additive fitness effects in section 6.3. This adapted model describes the fitness landscape well in section 6.4.

My understanding of the fitness landscape of synonymous mutations is, that the whole fitness landscape is dominated by the translation efficiency and within regions in the fitness landscape where the translation efficiency is constant, which are regions of equal travel time, the current is a secondary measure contributing to the fitness.

The description of the full fitness landscape from the paper by Zwart et al. [35] is this model's first success, but it needs to be further investigated in other, and hopefully larger, fitness landscapes, that my model accurately describes general fitness landscapes of synonymous mutations.

If the model proves to be true, either for systems that have bottlenecks far apart and far from the boundaries from chapter 5, or for systems that use the adaptation of additive fitness values within subcubes from section 6.3, this would give an interesting insight into the effect of synonymous mutations on the fitness landscape of translation.

Bibliography

- [1] R. J. Baker and R. D. Bradley, Speciation in Mammals and the Genetic Species Concept. *Journal of Mammalogy*, Volume 87, Issue 4, Pages 643–662 (2006), doi: <https://doi.org/10.1644/06-MAMM-F-038R2.1>
- [2] F. H. C. Crick, Central dogma of molecular biology. *Nature* 227, 561–563 (1970), doi: <https://doi.org/10.1038/227561a0>
- [3] J. F. Crow and W. F. Dove, Perspectives on Genetics: Anecdotal, Historical, and Critical Commentaries. *Univ of Wisconsin Press* (2000)
- [4] B. Derrida, E. Domany and D. Mukamel, An exact solution of a one-dimensional asymmetric exclusion model with open boundaries. *J. Stat. Phys.* 69, 667 (1992)
- [5] B. Derrida, M. R. Evans, V. Hakim and V. Pasquier, Exact solution of a 1d asymmetric exclusion model using a matrix formulation. *J. Phys. A, Math. Gen.* 26, 1493 (1993)
- [6] S. Foldes, A Characterization of Hypercubes. *Discrete Mathematics* Volume 17, Issue 2 1977, Pages 155–159 (1976)
- [7] M. E. Foulaadvand, S. Chaaboki and M. Saalehi Characteristics of the asymmetric simple exclusion process in the presence of quenched spatial disorder. *Phys. Rev. E* 75 (2007), doi: <https://doi.org/10.1103/PhysRevE.75.011127>
- [8] M. H. Friedman, Principles and Models of Biological Transport. *Springer, Berlin*, ISBN-10: 1441927158 (2008)
- [9] J. de Gier and F. H. L. Essler, Exact Spectral Gaps of the Asymmetric Exclusion Process with Open Boundaries. *Journal of Statistical Mechanics: Theory and Experiment* (2006) doi: <https://doi.org/10.1088/1742-5468/2006/12/p12011>
- [10] P. Greulich and A. Schadschneider, Phase diagram and edge effects in the ASEP with bottlenecks. *Physica A* 387: 1972-1986 (2008)
- [11] R. Hershberg and D. A. Petrov, Selection on Codon Bias. *Annual Review of Genetics* 42 1, 287-299 (2008), doi: <https://doi.org/10.1146/annurev.genet.42.110807.091442>
- [12] J. Howard, Mechanics of Motor Proteins and the Cytoskeleton. *Sinauer, Sunderland*, ISBN-10: 0878933336 (2001)
- [13] S. A. Janowsky and J. L. Lebowitz, Finite-size effects and shock fluctuations in the asymmetric simple-exclusion process. *Physical Review A* 45.2 (1992): 618.

- [14] S. A. Janowsky and J. L. Lebowitz, Exact results for the asymmetric simple exclusion process with a blockage. *Journal of Statistical Physics* 77, 35–51 (1994): 1572-9613
- [15] A. B. Kolomeisky and M. E. Fisher, Molecular motors:a theorist’s perspective. *Annu. Rev. Phys. Chem.* 58, 675 (2007)
- [16] J. Krug and L-H. Tang, Disorder-induced unbinding in confined geometries. *Phys. Rev. E* 50/1: 104-115, (1994), doi: <https://doi.org/10.1103/PhysRevE.50.104>
- [17] J. Krug, private communication (2019)
- [18] H. Liljenström and G. von Heijne Translation rate modification by preferential codon usage: Intragenic position effects. *Journal of Theoretical Biology* Volume 124, Issue 1, Pages 43-55, (1987), doi: [https://doi.org/10.1016/S0022-5193\(87\)80251-5](https://doi.org/10.1016/S0022-5193(87)80251-5)
- [19] C. T. MacDonald, J. H. Gibbs and A. C. Pipkin, Kinetics of biopolymerization on nucleic acid templates. *Biopolymers* 6(1) (1968)
- [20] C.T. MacDonald and J. H. Gibbs, Concerning the kinetics of polypeptide synthesis on polyribosomes. *Biopolymers* 7, 707 (1969)
- [21] J. L. Parmley and L. D. Hurst, How do synonymous mutations affect fitness? *Bioessays*, 29: 515-519 (2007), doi: <https://doi.org/10.1002/bies.20592>
- [22] J. B. Plotkin and G. Kudla, Synonymous but not the same: the causes and consequences of codon bias. *Nature Reviews Genetics* 12, 32, (2010), doi: <https://doi.org/10.1038/nrg2899>
- [23] A. Schadschneider, D. Chowdhury and K. Nishinari, Stochastic Transport in Complex Systems. *Elsevier*, ISBN-13: 9780444528537 (2011)
- [24] M. F. Schenk, I. G. Szendro, J. Krug and J. A. G. M. de Visser, Quantifying the Adaptive Potential of an Antibiotic Resistance Enzyme. *PLOS Genetics* 8(6): e1002783. (2012), doi: <https://doi.org/10.1371/journal.pgen.1002783>
- [25] M. F. Schenk, I. G. Szendro, M. L. M. Salverda, J. Krug and J. A. G. M. de Visser, Patterns of Epistasis between Beneficial Mutations in an Antibiotic Resistance Gene. *Molecular Biology and Evolution*, Volume 30, Issue 8, August 2013, Pages 1779–1787 (2013), doi: <https://doi.org/10.1093/molbev/mst096>
- [26] M. Schliwa and G. Woehlke, Molecular motors. *Nature* 422, 759 (2003)

Bibliography

- [27] J. Szavits-Nossan, Disordered exclusion process revisited: some exact results in the low-current regime. *Journal of Physics A: Mathematical and Theoretical* 46.31 (2013): 315001
- [28] J. Szavits-Nossan, M. C. Romano and L. Ciandrini, Power series solution of the inhomogeneous exclusion process. *Phys. Rev. E* 97, 052139 (2018), doi: <https://doi.org/10.1103/PhysRevE.97.052139>
- [29] J. Szavits-Nossan and M. R. Evans, Dynamics of ribosomes in mRNA translation under steady- and nonsteady-state conditions. *Phys. Rev. E*, volume 101,6 *American Physical Society* (2020), doi: <https://doi.org/10.1103/PhysRevE.101.062404>,
- [30] P. Taylor and R. Lewontin, The Genotype/Phenotype Distinction. *The Stanford Encyclopedia of Philosophy, Metaphysics Research Lab, Stanford University* (Summer 2017), doi: <https://plato.stanford.edu/archives/sum2017/entries/genotype-phenotype/>
- [31] S. Varenne and C. Lazdunski, Effect of distribution of unfavourable codons on the maximum rate of gene expression by an heterologous organism. *Journal of Theoretical Biology* Volume 120, Issue 1, Pages 99-110 (1986), doi: [https://doi.org/10.1016/S0022-5193\(86\)80020-0](https://doi.org/10.1016/S0022-5193(86)80020-0)
- [32] D. M. Weinreich, R. A. Watson and L. Chao, Perspective: sign epistasis and genetic constraint on evolutionary trajectories. *Evolution* 59.6 (2005): 1165-1174
- [33] S. Wright, The roles of mutation, inbreeding, crossbreeding, and selection in evolution. p. 356-366 (1932)
- [34] R. K. P. Zia, J. J. Dong and B. Schmittmann, Modeling Translation in Protein Synthesis with TASEP: A Tutorial and Recent Developments. *J Stat Phys* 144: 405 (2011), doi: <https://doi.org/10.1007/s10955-011-0183-1>
- [35] M. P. Zwart, M. F. Schenk, S. Hwang et al., Unraveling the causes of adaptive benefits of synonymous mutations in TEM-1 β -lactamase. *Heredity* 121, 406–421 (2018), doi: <https://doi.org/10.1038/s41437-018-0104-z>

A. Appendix: Description of code simulating the TASEP

The code to simulate the systems discussed in this thesis was written using python. I use this code to both simulate already published examples to prove the code's validity and also get a general understanding of the details of a TASEP with multiple bottlenecks. Finally the results from these simulations are used to validate a model of interacting slow sites in a TASEP and compare it to experimental results.

The traditional way of simulation a TASEP is performed using Monte Carlo simulations. Each site i is associated with a jumping probability ω_i . Within one Monte Carlo step each site will be picked on average once and it is tested if

- a) there is a particle at that site i ($\kappa_i = 1$), with κ_i the occupancy of site i ,
- b) there is no particle at the next site ($\kappa_{i+1} = 0$) and
- c) a randomly generated number is smaller than the jumping probability , with κ_i the occupancy of site i

If all of these conditions are true, the particle jumps to the next site. This means that the probability of jumping with every try is

$$P_{\text{jump},i} = \rho_i(1 - \rho_{i+1})r_i , \quad (\text{A.1})$$

with the densities ρ_i, ρ_{i+1} at sites $i, i + 1$ and the rate r_i at site i .

This method is quite slow when the system is far from the maximum current phase, as there are few particles ready to jump available leading to long waiting times between updates. Therefore I used a different algorithm.

My program keeps track of all sites that are potential jumpers while only calculating successful jumps. Every step of the program begins by selecting randomly among those sites that could jump (conditions **a**) and **b**) are met) weighted by r_i , r_i is the jump rate at site i . Instead of checking if the conditions above apply a in every time step, i calculate the average time that the system would have taken until the next successful jump and add this to the list of times before selecting the next successful jumper randomly weighted by the jump rates of each potential jumper. Condition **c**) is no longer required as I use the expected value of the waiting time instead of a random number to have the system take another step.

The only random element in this simulation is which potential jumper is the next to move.

NEW NAME A. APPENDIX: DESCRIPTION OF CODE SIMULATING THE TASEP

This replaces the waiting time distribution due to the randomness of selecting a suitable site by its average over all time-steps. As we are interested in long times, averaging the time until the next step should not change the result.

The waiting time, i.e. the time that the system is expected to do nothing between successful jumps, averaged over all jumps can be calculated with

$$\tau = \left\langle \frac{1}{\sum_i r_i \delta_{\kappa_i,1} \delta_{\kappa_{i+1},0}} \right\rangle_t, \quad (\text{A.2})$$

with the Kronecker deltas conditioning the sum to only count the rates where conditions **a)** and **b)** meet and t the number of jumps.

The equilibrium current J of a system of length L is

$$J = \frac{1}{\tau L}. \quad (\text{A.3})$$

The density is calculated directly by averaging the occupancy κ_i of the sites over all successful jumps

$$\rho_i = \langle \kappa_i \rangle \quad (\text{A.4})$$

and the average density in the system is

$$\bar{\rho} = \frac{\sum_i \kappa_i}{L}. \quad (\text{A.5})$$

B. Appendix: Results of the search algorithm for different densities after the bottlenecks

In section 4.4 the search algorithm to find interfaces within the phase space is described. The following results are different runs of the search algorithm, which show more phase space examples for different density values. The graphs below are calculated with my search algorithm, that searches the phase space for point quartets, where some of the four points in the phase space have different orderings of the travel times $\bar{\tau}_0, \bar{\tau}_1, \bar{\tau}_2$ than the others. If all four points have the same ordering, the algorithm continues traveling left, if it started on the right, and to the right if it started on the left. The upper part ($\rho_2 > x_2$) is the one that is investigated more thoroughly because at the time of calculation, the goal was to find the area in the phase space where $\bar{\tau}_2 < \bar{\tau}_1 < \bar{\tau}_0$.

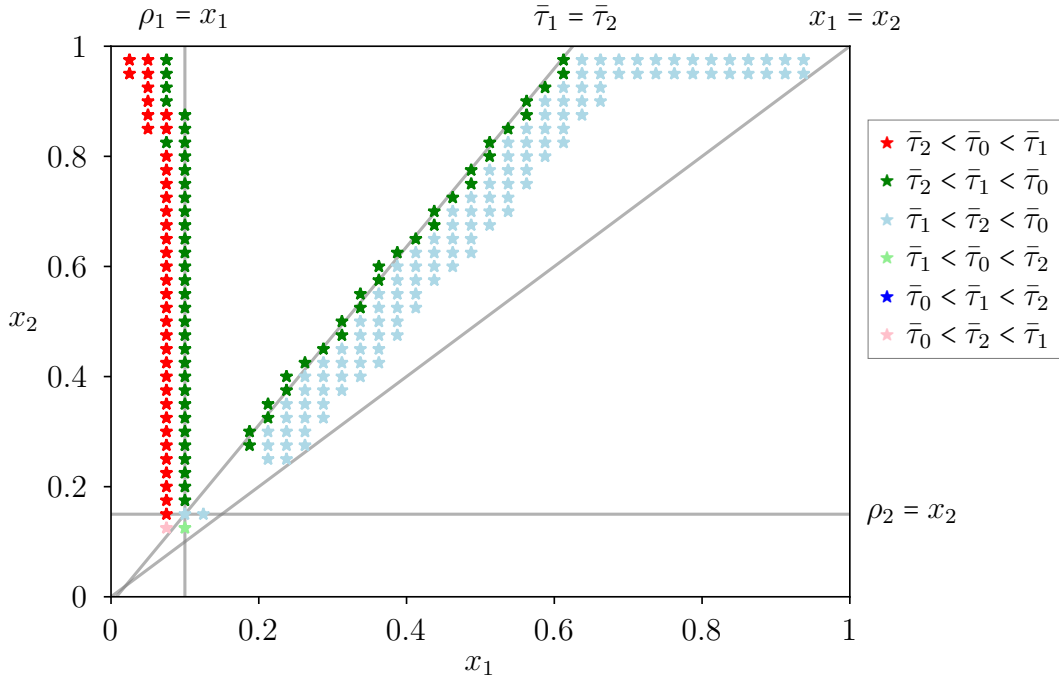


Figure B.1.: The phase space for parameters $\rho_1 = 0.1$ and $\rho_2 = 0.15$ with system size $L = 800$. The result shows the line $\rho_1 = x_1$ and $\bar{\tau}_1 = \bar{\tau}_2$ nicely, but does not show the lower values, because the points in the landscape were too close together.

NEW NAME B. APPENDIX: RESULTS OF THE SEARCH ALGORITHM FOR DIFFERENT DENSITIES AFTER THE BOTTLENECKS

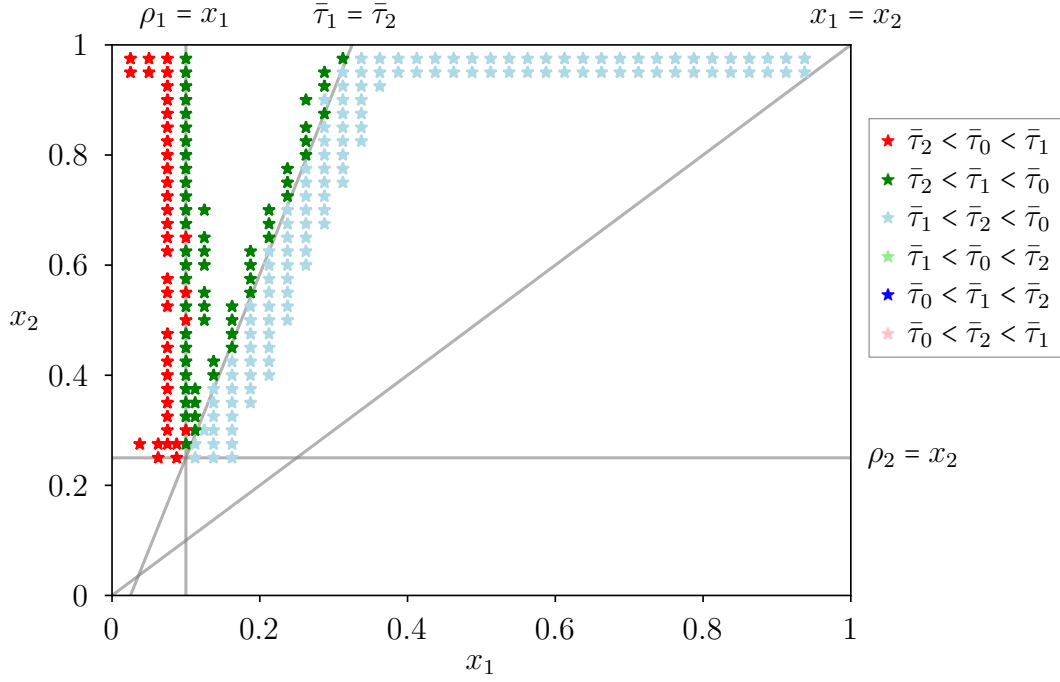


Figure B.2.: The phase space for parameters $\rho_1 = 0.1$ and $\rho_2 = 0.25$ with system size $L = 800$. This version of the search algorithm was not prepared for finding the phase transitions for $\rho_2 < x_2$. The bottom part of the lines can only be shown with an updated version of the algorithm which was not written at the time.

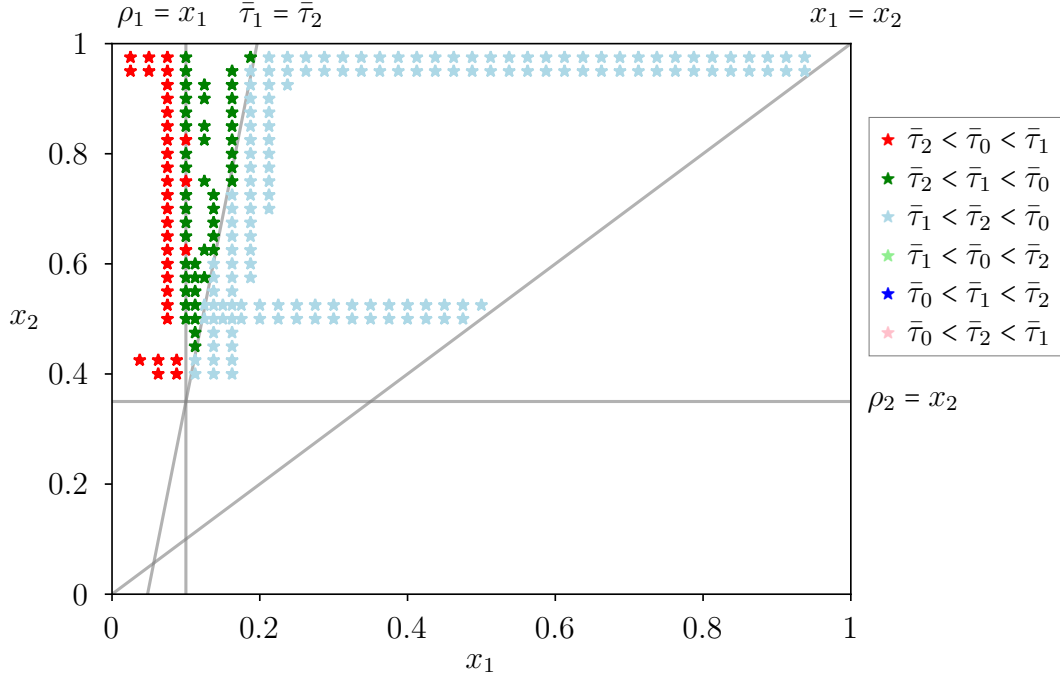


Figure B.3.: The phase space for parameters $\rho_1 = 0.2$ and $\rho_2 = 0.25$ with system size $L = 800$. As above, the three lines intersect and this version of the algorithm stops there, even after a restart closer to the intersection (second row of red and light-blue points).

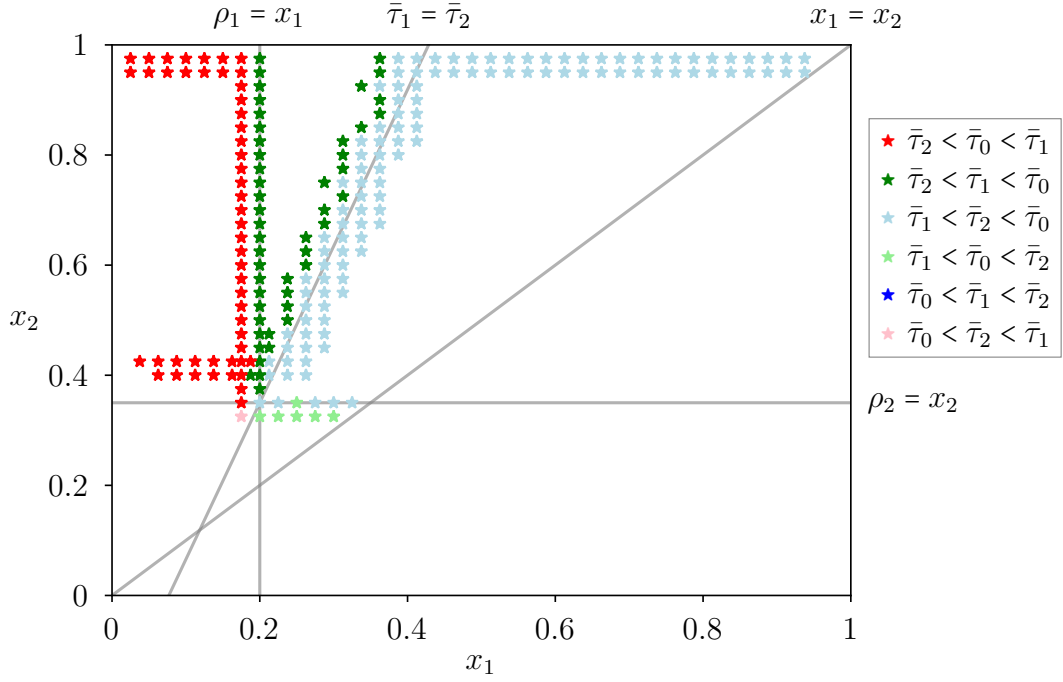


Figure B.4.: The phase space for parameters $\rho_1 = 0.2$ and $\rho_2 = 0.35$ with system size $L = 800$. The second start from the right shows the line $\rho_2 = x_2$.

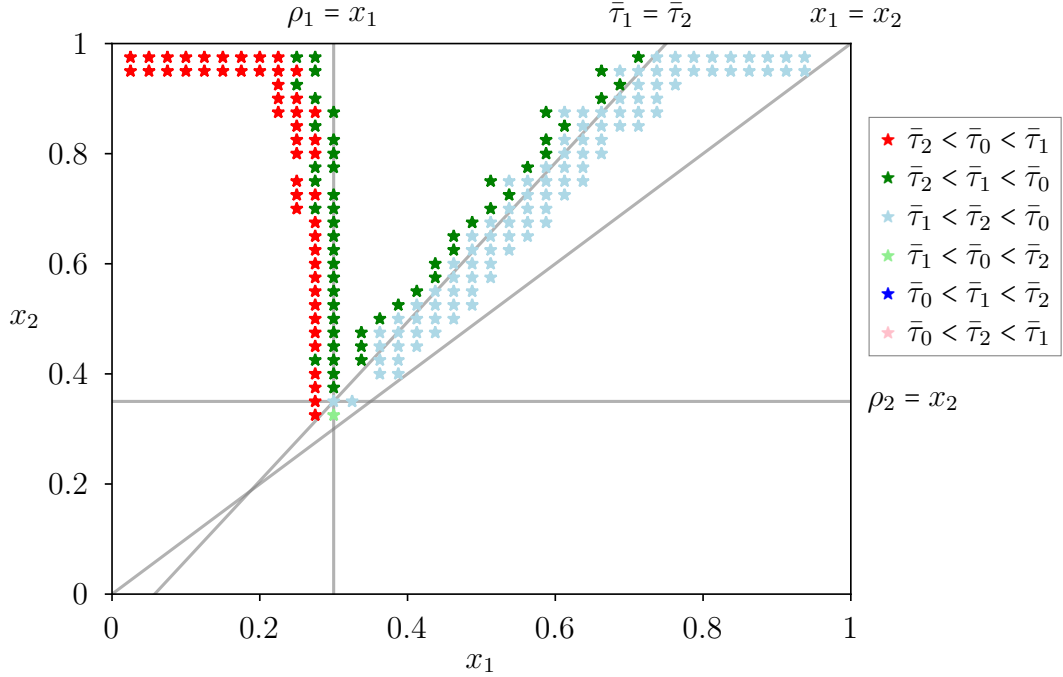


Figure B.5.: The phase space for parameters $\rho_1 = 0.3$ and $\rho_2 = 0.35$ with system size $L = 800$. Here the intersection is too close to $x_1 = x_2$ to show more than the lines $\rho_1 = x_1$ and $\bar{\tau}_1 = \bar{\tau}_2$.

Acknowledgment

I want to thank my whole group for the discussions we had. Both the scientific and also the clearly non-scientific discussions were really memorable and enjoyable. Thank you for listening to me when I tried to explain my ideas on my topic and for listening to me ramble about cooking or other private matters.

I would like to thank to my office colleagues Lucy, Alex and Till, for making my time, when it was still possible to go to the office, both fun and productive.

I want to thank my supervisor Joachim for giving me the opportunity to write on this nice topic, for always having an open door and ear for my questions and for his incredible insight, of which I tried to absorb as much as I could.

Finally I want to thank my girlfriend Eva and my friends Stefan and Shayoni for checking my thesis for mistakes and helping me understand, that some things are not as trivial as I imagined.